



Exceptional service in the national interest

The Promise of Neuromorphic Computing

Rob Hoekstra

June 26, 2023



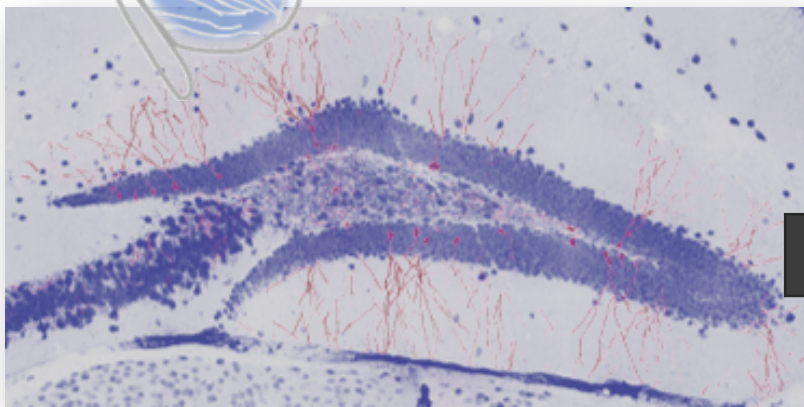
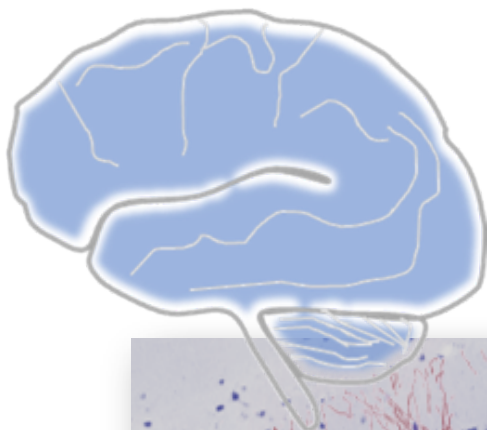
So why the brain?



- Energy efficient
- Operationally fast considering slow components
- Data efficient
- Diverse applications
- Robustness

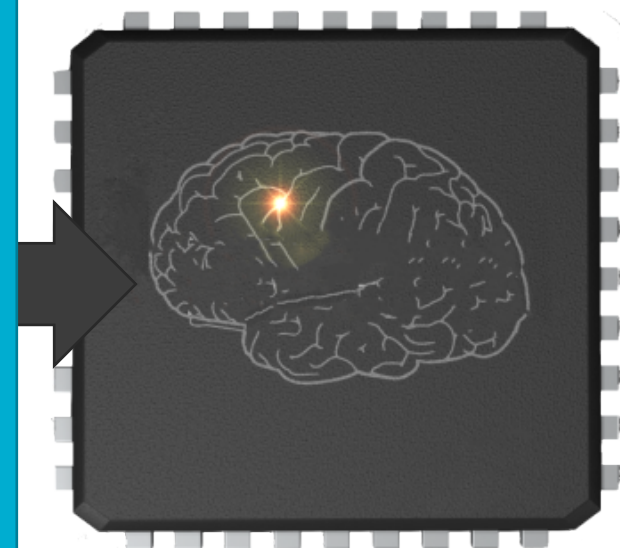


What is the brain good at?



Realized Features of Brain Inspiration in Neuromorphic Hardware

- Event-driven communication
- Graph based connectivity
- Processing in Memory
- In situ learning
- Analog computation
- Post-Moore's Law Devices
- Ubiquitous stochasticity



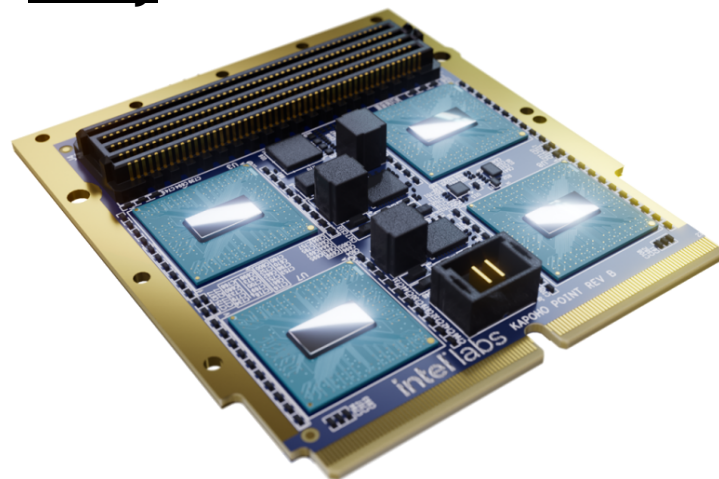


Where are we today in neuromorphic computing?

Realized Features of Brain Inspiration in Neuromorphic Hardware

- Event-driven communication
- Graph based connectivity
- Processing in Memory
- In situ learning
- Analog computation
- Post-Moore's Law Devices
- Ubiquitous stochasticity

Today

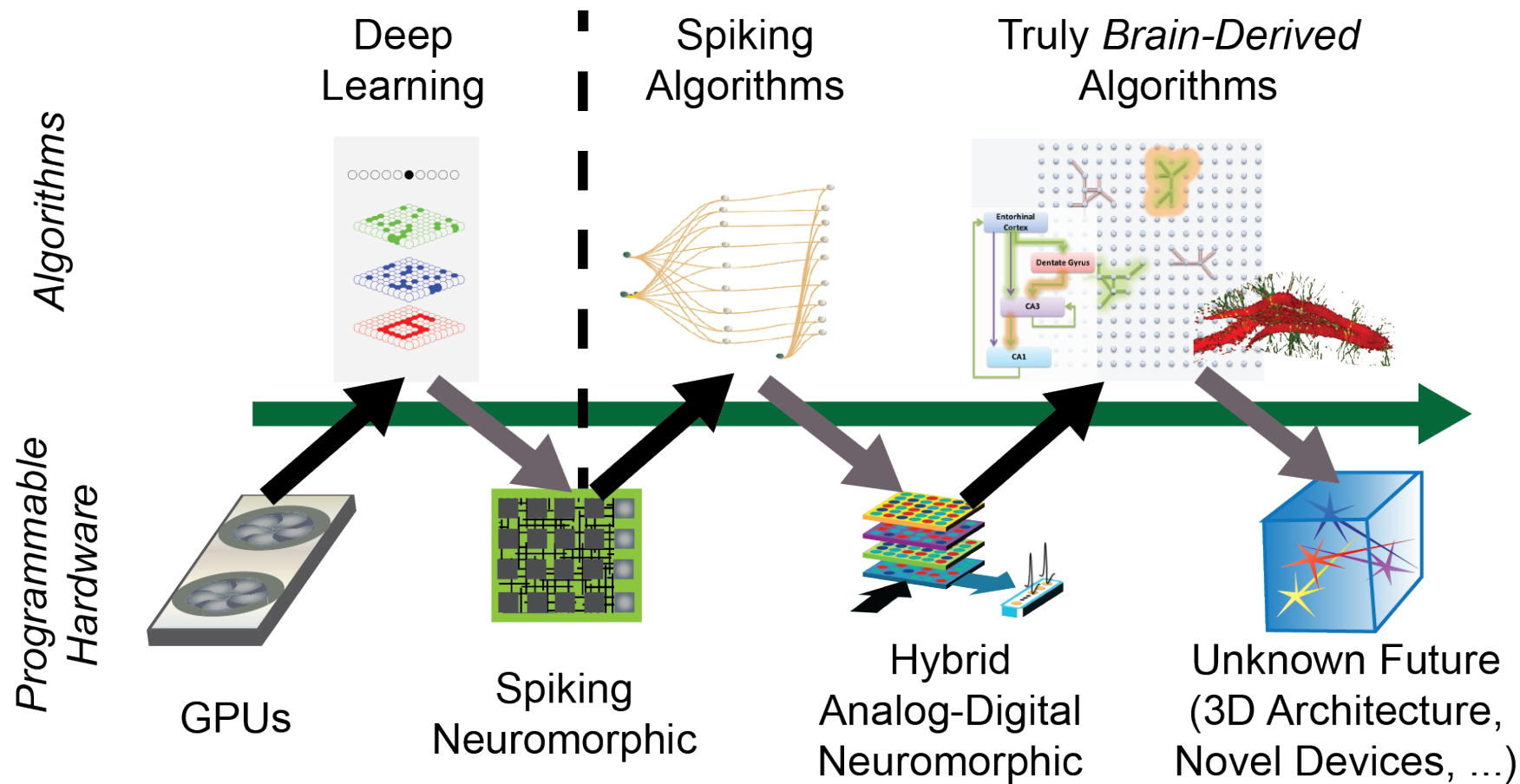


Intel Loihi 2
Millions of CMOS neurons
Billions of CMOS synapses
~ 1 Watt power

One example (of many...)

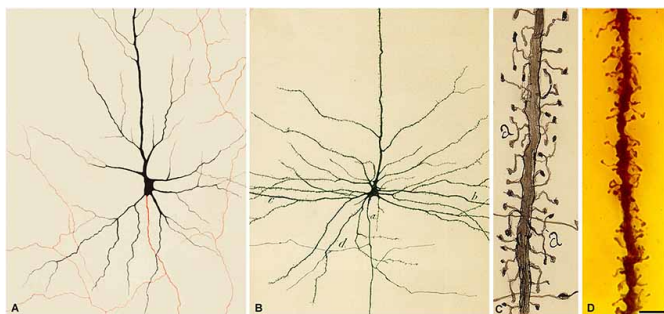
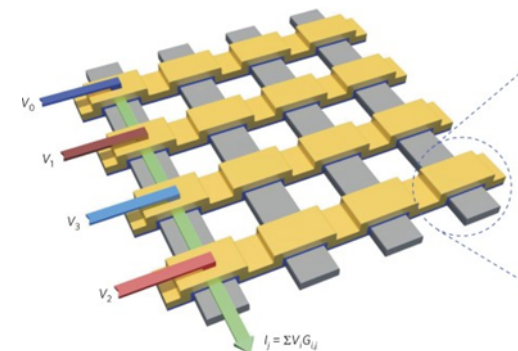
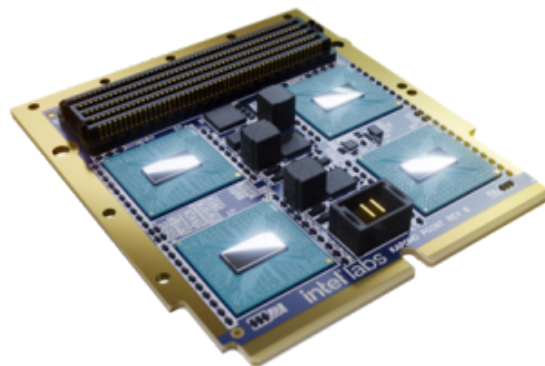
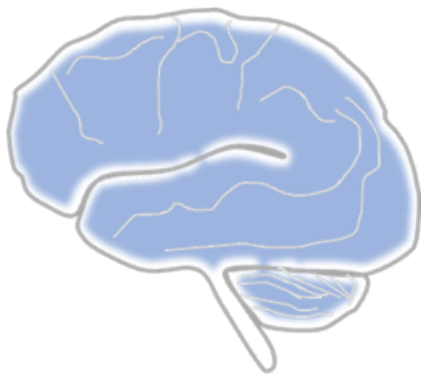


Towards novel future architectures

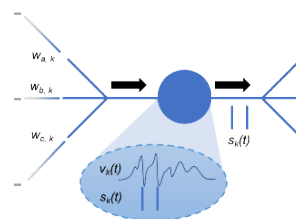




Fundamental science question: what is the scale of neuromorphic computing needed?



$>10^{11}$ neurons
 $>10^{15}$ synapses
High complexity
Highly efficient
Slow



$>10^5$ neurons
 $>10^8$ synapses
Low complexity
Moderately efficient
Fast

$$I = \Sigma V/R$$

$\sim 10^2$ neurons
 $\sim 10^4$ synapses
Very low complexity
Highly efficient
Fast?

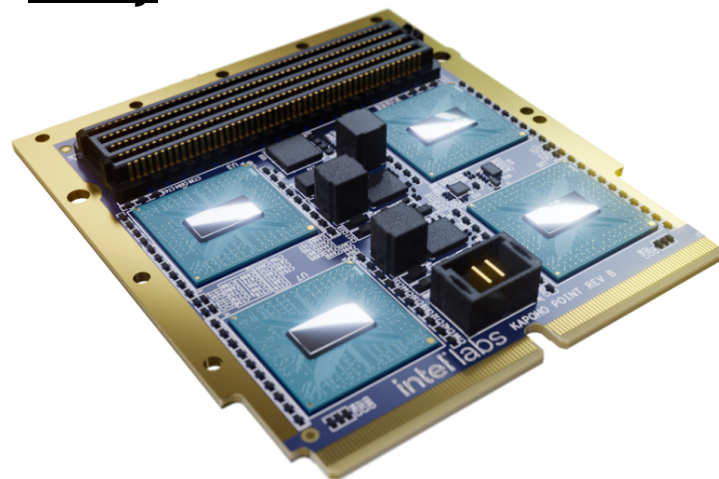


Where are we today in neuromorphic computing? And where might we be in the future?

Realized Features of Brain Inspiration in Neuromorphic Hardware

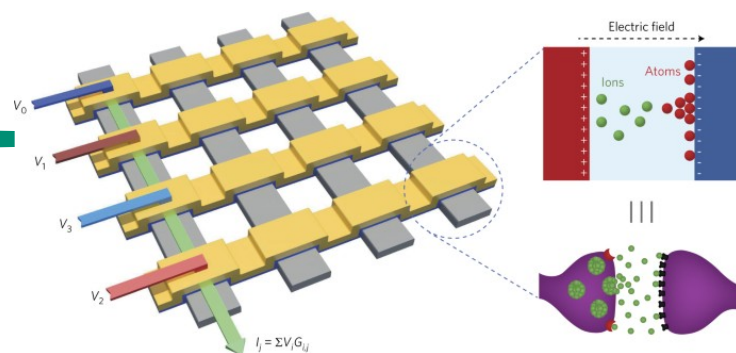
- Event-driven communication
- Graph based connectivity
- Processing in Memory
- In situ learning
- Analog computation
- Post-Moore's Law Devices
- Ubiquitous stochasticity

Today



Intel Loihi 2
Millions of CMOS neurons
Billions of CMOS synapses
~ 1 Watt power

Tomorrow



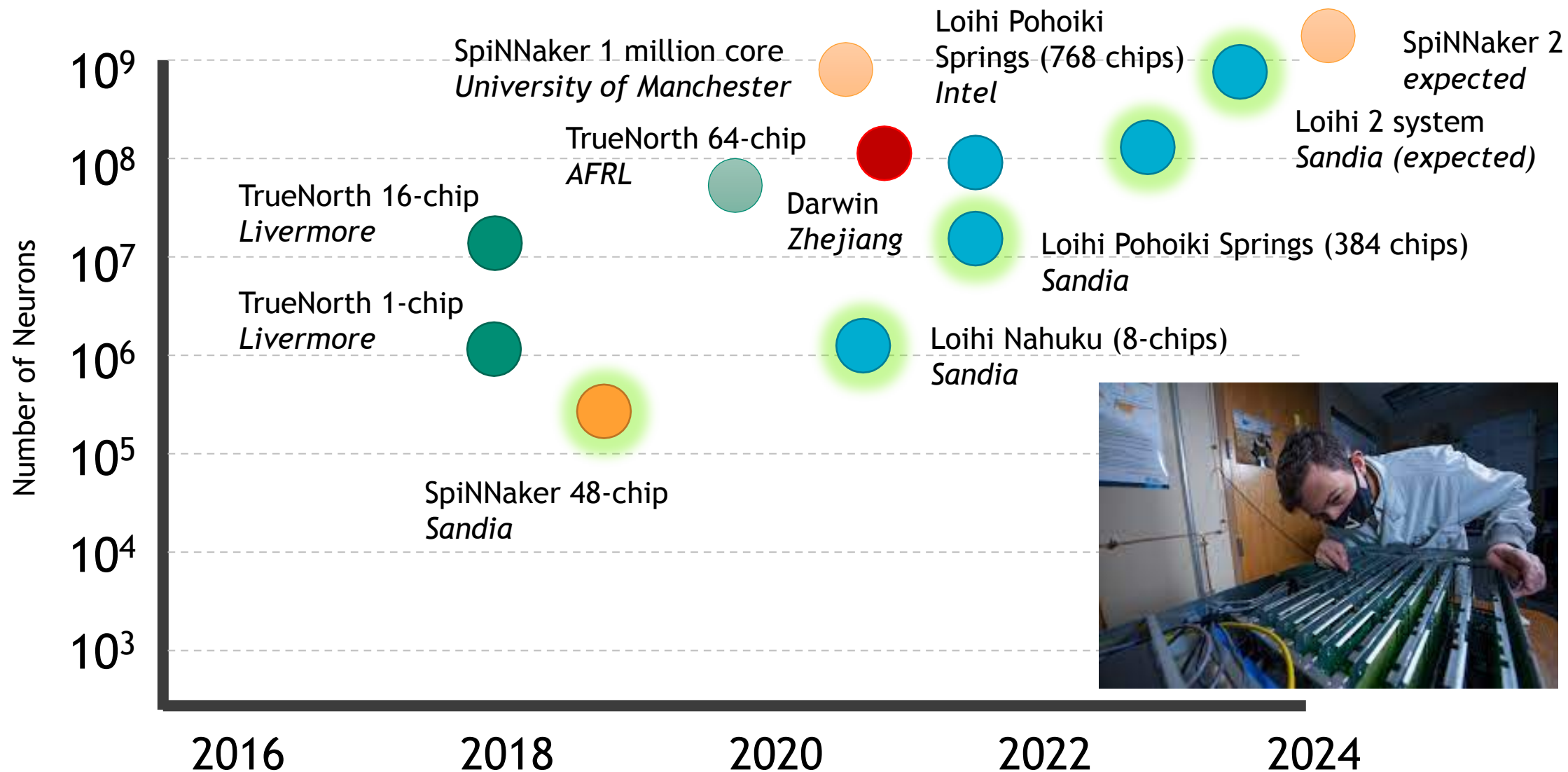
Post-Moore Devices
(ECRAM, Memristors, MTJs,
optical, organic, etc)

Scale to human sizes?

Zidan et al., 2018

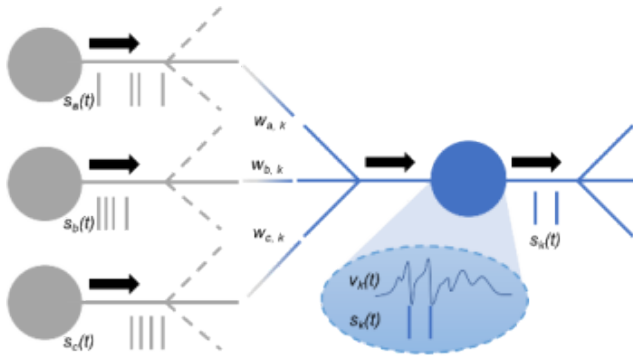


Sandia hosts some of the world's largest CMOS-based neuromorphic systems





Spiking neuromorphic today: Overview

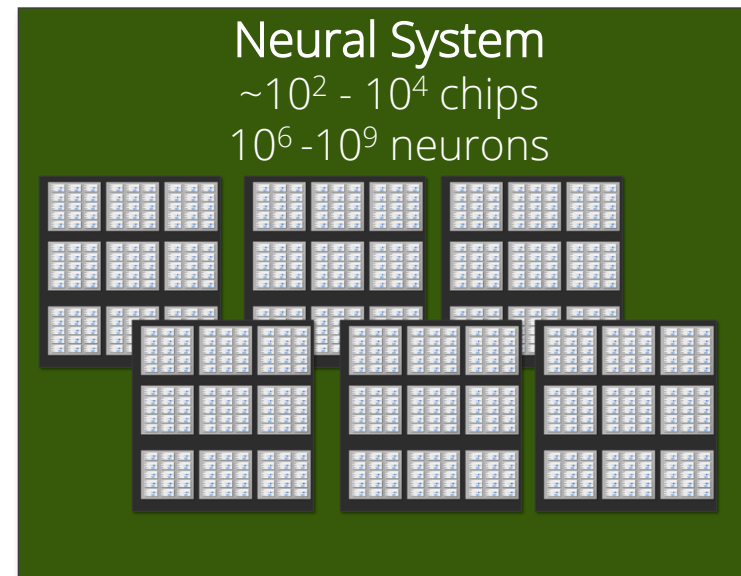
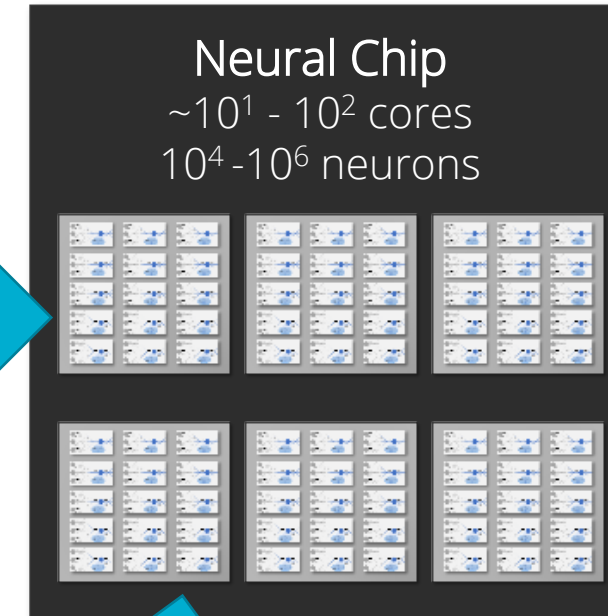
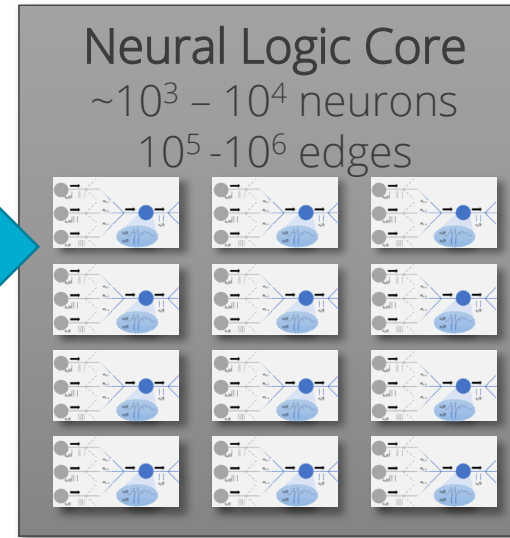


Computational Primitives:

Spiking Neurons (vertices / nodes)
Synapses (connections / edges)

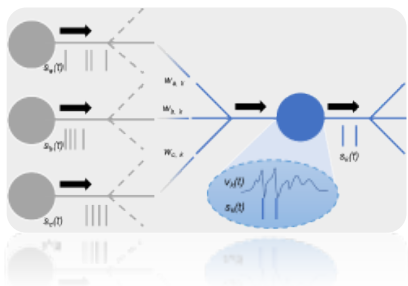
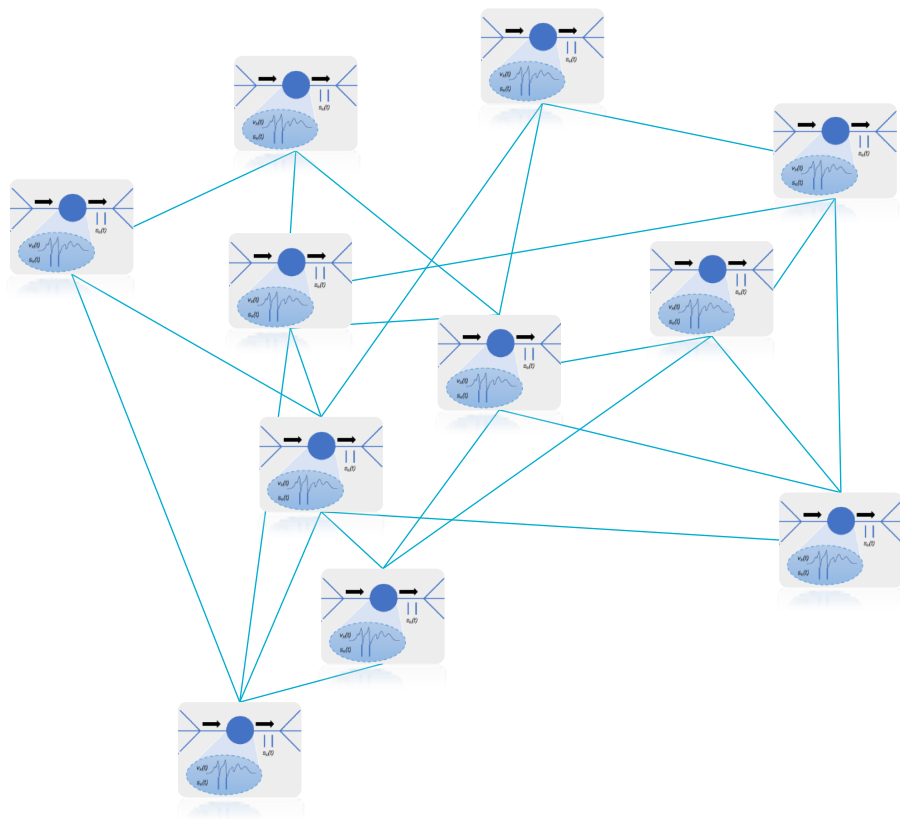
Programmable as arbitrary graphs

- Edges: Directed and weighted
- Nodes: Threshold gate logic + time
- *Artificial neural networks are a special case*
- Programmability, theoretical, analysis and software are open research questions





Neuromorphic hardware jumped ahead of the rest of the stack



Neuromorphic hardware has been built with a “if we build it, neuroscientists will come” hope

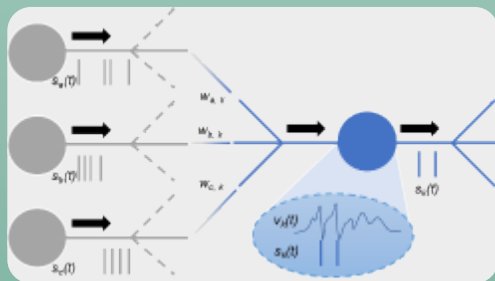
We need

- ❖ Driving Applications
- ❖ Systems Interface
- ❖ Software and Programming Paradigm
- ❖ Theoretical Framework

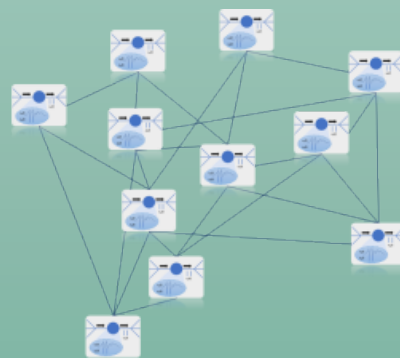


A quick aside: most neuromorphic hardware is **not** designed for current artificial neural networks

Neuromorphic Hardware



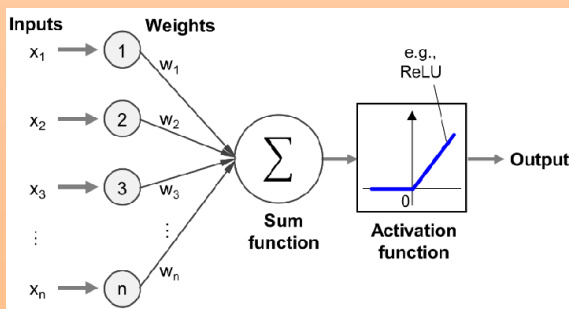
Spiking neurons



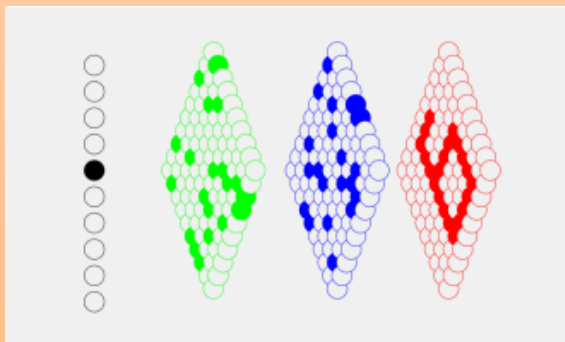
Arbitrary connectivity

- Continual learning integrated into operation
- Inherently temporal
- Dynamical tasks?

Artificial Neural Networks



Continuous neurons



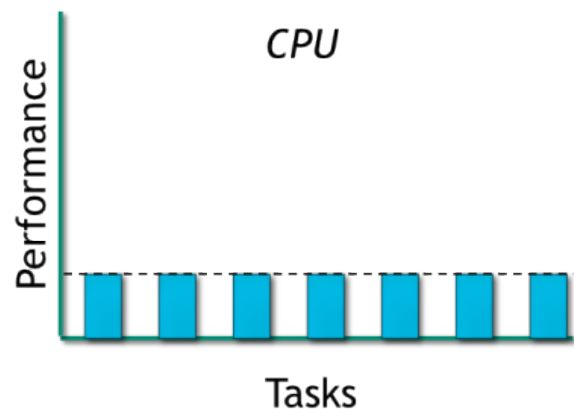
Linear algebra-like networks

- Distinct training and inference modes
- Time is largely avoided
- Computer vision and natural language

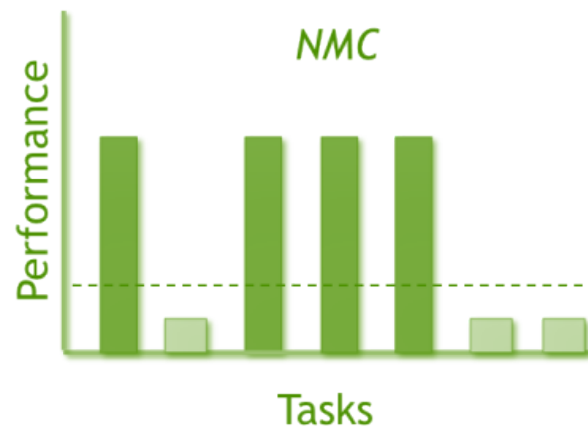
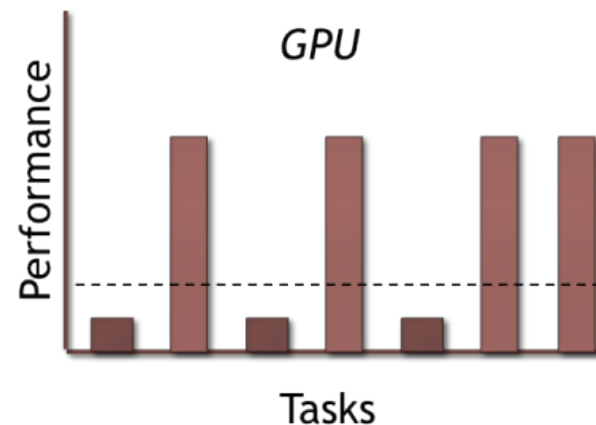


Could neuromorphic be generalized to more algorithms just as GPUs have been?

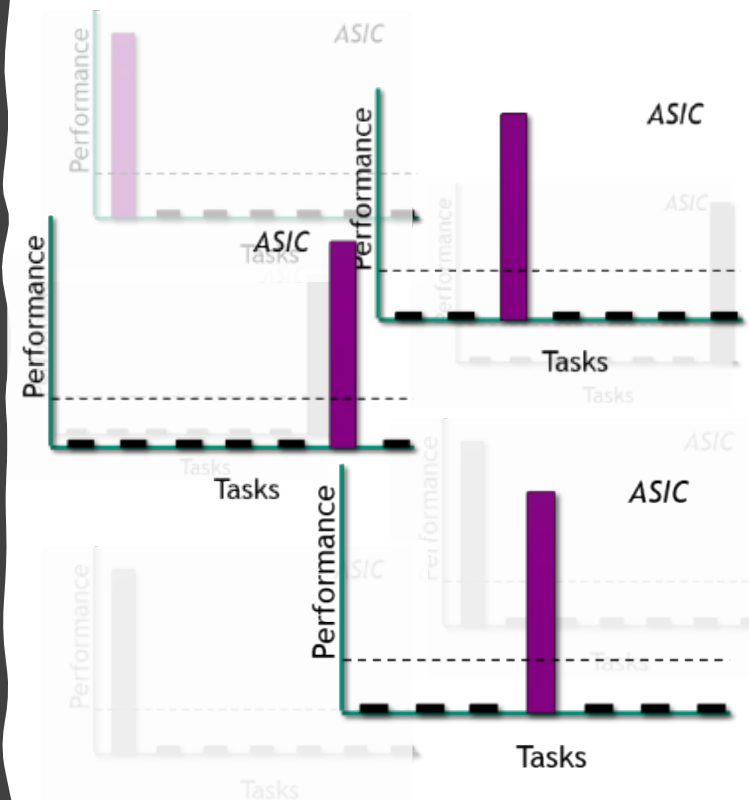
Truly General Purpose



Specialized General Purpose



Application Specific

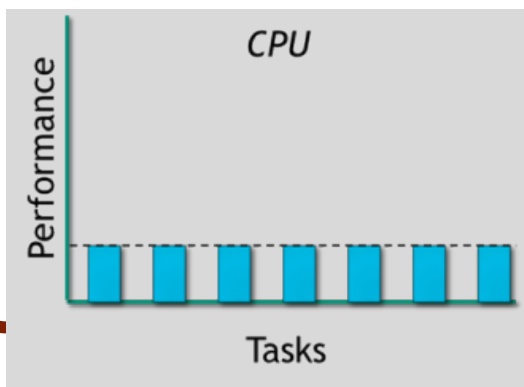




Separating the “can do” from the “should do”

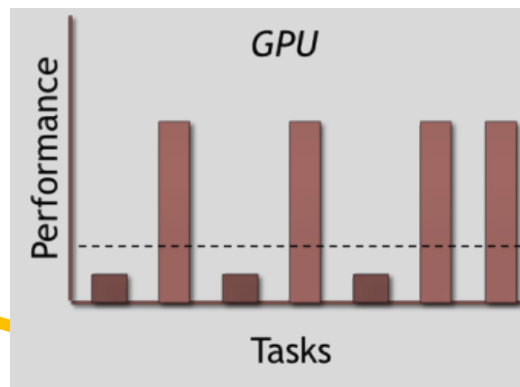
Can implement on NMC, but only to avoid I/O

- Arithmetic (adding, subtraction, multiplication, etc.)
- Data filtering
- Sorting
- Data conversions
- ...



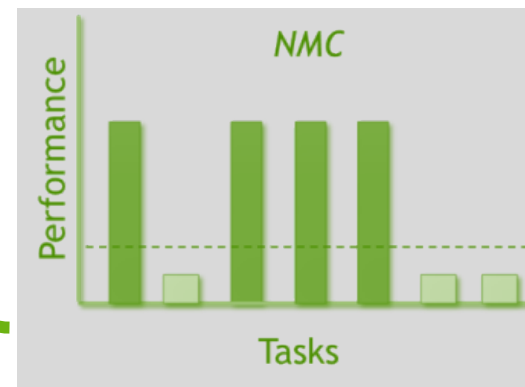
Possibly good on NMC, but there may be alternatives

- Deep learning / conventional artificial neural networks
- Parallel data processing (background and change detection, convolutions, etc)
- Linear algebra (MVM, cross-correlations, L1-norm, etc)
- Classic machine learning (SVMs, k-nearest neighbors, clustering)



Should implement on NMC once systems reach scale

- Algorithms the brain actually uses (** we don't have these yet...*)
- Random walks / Discrete Time Monte Carlo
- *Some* Graph Algorithms (Dynamic programming, Dijkstra, triangle counting, graph cut, etc)
- *Some* neural networks





Neuromorphic computing can impact a broad range of applications

 IOPscience

Neuromorphic Computing and Engineering

ACCEPTED MANUSCRIPT • OPEN ACCESS

A review of non-cognitive applications for neuromorphic computing

James Aimone¹ , Prasanna Date², Gabriel Fonseca-Guerra³, Kathleen Hamilton², Kyle Henke⁴, Bill Kay⁵, Garrett Kenyon⁴, Shruti Kulkarni², Susan Mniszewski⁶, Maryam Parsa⁷, Sumedh Risbud³ , Catherine Schuman⁸ , William Severa¹ and J. Darby Smith¹ [— Hide full author list](#)

Accepted Manuscript online 10 August 2022 • © 2022 The Author(s). Published by IOP Publishing Ltd

64 Total downloads



[Turn on MathJax](#)

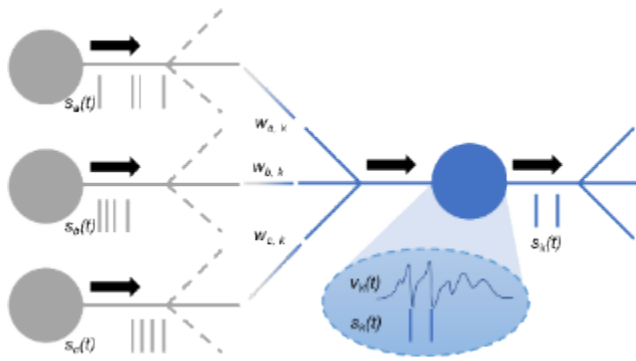
Share this article



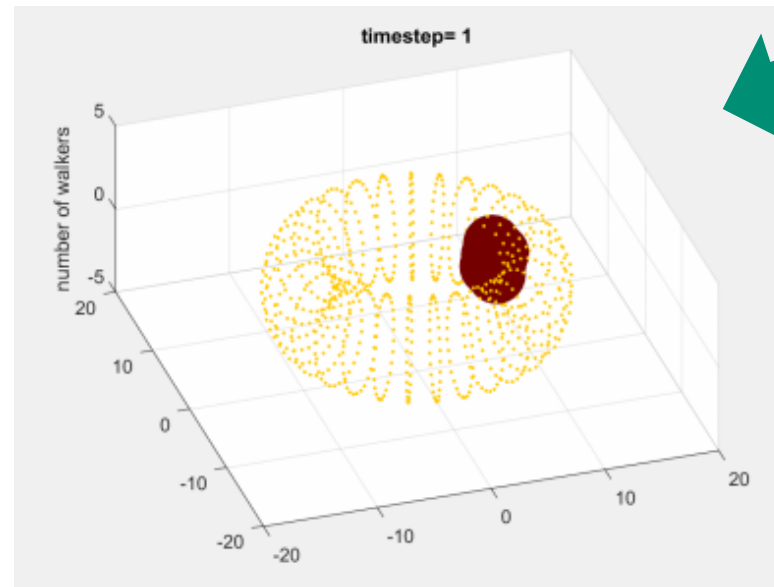
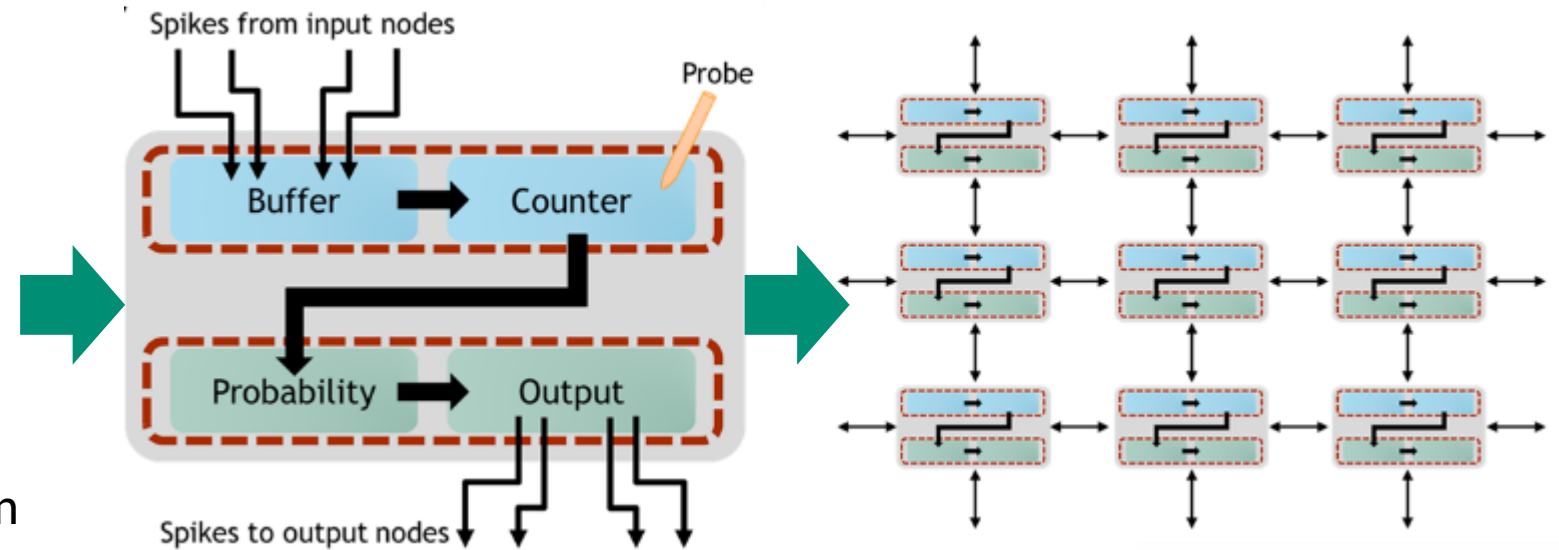
*Spiking
Scientific
Computing*



Today's spiking NMC shows energy advantage over conventional approaches on Monte Carlo simulations



Leaky Integrate and Fire Neuron



Neuromorphic scaling advantages for energy-efficient random walk computations

J. Darby Smith, Aaron J. Hill, Leah E. Reeder, Brian C. Franke, Richard B. Lehoucq, Ojas Parekh, William Severa and James B. Alimonte

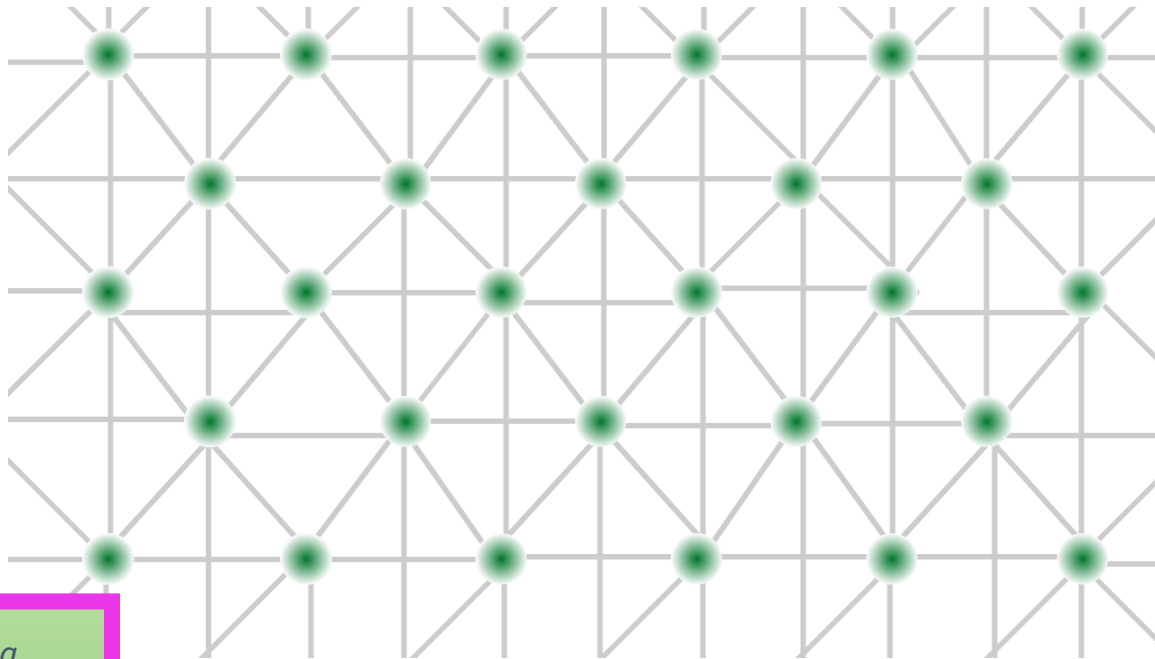
Neuromorphic computing, which aims to replicate the computational structure and architecture of the brain in synthetic hardware, has typically focused on artificial intelligence applications. What is less explored is whether such brain-inspired hardware can provide value beyond cognitive tasks. Here we show that the high degree of parallelism and configurability of spiking neuromorphic architectures makes them well suited to implement random walks via discrete-time Markov chains. These random walks are useful in Monte Carlo methods, which represent a fundamental computational tool for solving a wide range of numerical computing tasks. Using IBM's TrueNorth and Intel's Loihi neuromorphic computing platforms, we show that our neuromorphic computing algorithm for generating random walk approximations of diffusion offers advantages in energy-efficient computation compared with conventional approaches. We also show that our neuromorphic computing algorithm can be extended to more sophisticated jump-diffusion processes that are useful in a range of applications, including financial economics, particle physics and machine learning.

Despite the increasing ability to develop large-scale neural processors today¹, the theoretical value of neuromorphic hardware remains unclear—unlike quantum computing that offers clear fundamental advantages at scale². Nevertheless, there are several architectural features of most nervous systems that could yield advantages including the high degree of connectivity between neurons, the collocation of processing and memory, and the use of action potentials (referred to as spikes) to communicate^{3,4}. Algorithm research for spiking neuromorphic hardware has primarily focused on its suitability for deep learning and other emerging artificial intelligence (AI) algorithms^{5–7}. Such applications are straightforward, given the alignment of neural architectures with neural networks, and it can be expected that the value of neuromorphic computing will grow as AI algorithms derive further inspiration from the brain⁸. However, the impact of neuromorphic computing beyond cognitive applications is less certain. Quantum computing has shown how emerging hardware can have an impact beyond its original inspiration: it was conceived as a means for efficient chemistry simulations^{9,10}, but is now recognised as useful in a much broader range of applications^{11–13}. Unlike quantum computing, which faces technical challenges in scaling up¹⁴, time scaling compared with the von Neumann architecture and still requiring less total energy to perform the same computation. Observing a neuromorphic advantage for non-cognitive applications should not be taken as a given since the specialization of computer architectures to improve performance on a subset of tasks will likely result in degraded performance in other tasks¹⁵. Therefore, observing a neuromorphic advantage on non-cognitive applications would demonstrate that neuromorphic computing can have a broader impact than previously assumed and provide a concrete framework by which to develop the technology. Although a definitive neuromorphic advantage (as defined here) has not yet been demonstrated for non-cognitive applications, there are three categories of such computing tasks that appear well suited for neuromorphic computing: linear algebra, in which the high fan-in of neurons can be used to realise known theoretical advantages of threshold gate (TG) logic^{16,17}; graph analytical tasks that can leverage the configurability and parallelization of neural circuits^{18–21}; and sampling steady-state distributions for a wide range of potential applications using stochastic neural circuits^{22–24}. In this Article, we show that large-scale neuromorphic hardware can offer a neuromorphic advantage on a fundamental

Spiking
Scientific
Computing



Neuromorphic computing advantage appears to be when an algorithm can split task across computational graph with sparse communication



Spiking
Scientific
Computing

Monte Carlo simulations

Discrete Time Markov Chains

Dynamic programming

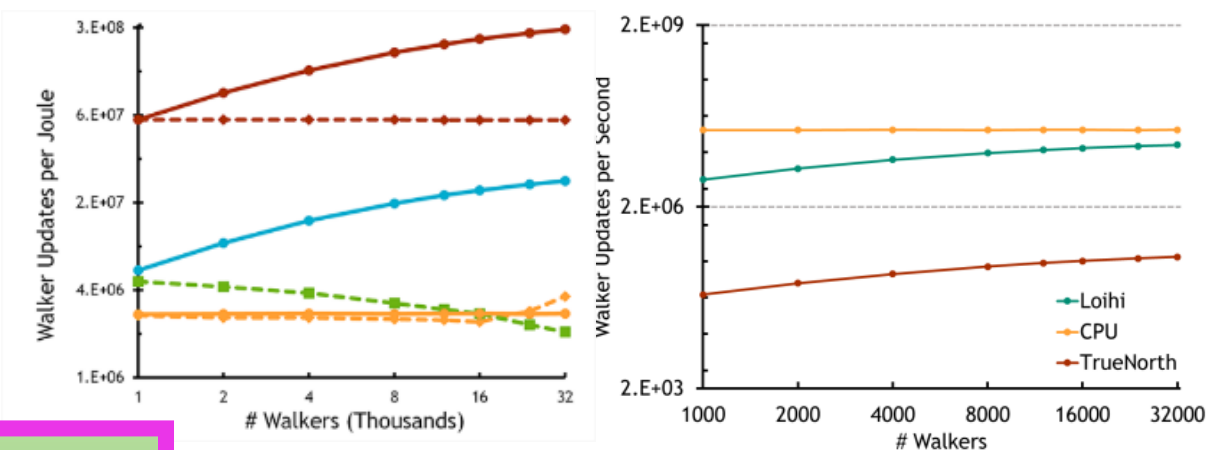
Graph neural networks

...

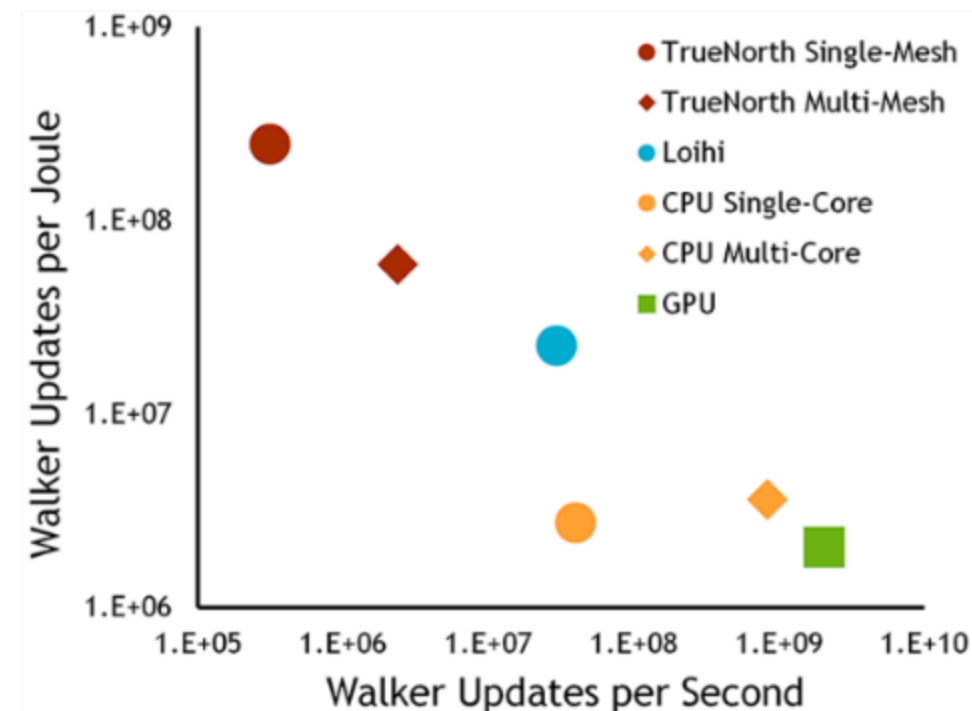


We can identify a neuromorphic advantage for simulating random walks

We define a *neuromorphic advantage* as an algorithm that shows a demonstrable **advantage** in terms of one resource (e.g., energy) while exhibiting comparable **scaling** in other resources (e.g., time).



Spiking
Scientific
Computing





Math: What PDEs can these stochastic processes be useful for?

Class of Partial Integro-Differential Equations:

$$\begin{aligned} \frac{\partial}{\partial t} u(t, \mathbf{x}) = & \frac{1}{2} \sum_{i,j} (\mathbf{a}\mathbf{a}^\top)_{i,j}(t, \mathbf{x}) \frac{\partial^2}{\partial x_i \partial x_j} u(t, \mathbf{x}) + \sum_i b_i(t, \mathbf{x}) \frac{\partial}{\partial x_i} u(t, \mathbf{x}) \\ & + \lambda(t, \mathbf{x}) u(t, \mathbf{x}) + \int \mathbf{h}(t, \mathbf{x}, q) [u(t, \mathbf{x} + q) - u(t, \mathbf{x})] \phi_Q(q; t, \mathbf{x}) dq \\ & + c(t, \mathbf{x}) u(t, \mathbf{x}) + f(t, \mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d, t \in [0, \infty). \end{aligned}$$

Stochastic Process:

NMC Hardware Simulates This Stochastic Process

$$d\mathbf{X}(t) = \mathbf{b}(t, \mathbf{X}(t)) dt + \mathbf{a}(t, \mathbf{X}(t)) d\mathbf{W}(t) + \int \mathbf{h}(t, \mathbf{X}(t), q) dP(t; Q, \mathbf{X}(t)).$$

Spiking
Scientific
Computing

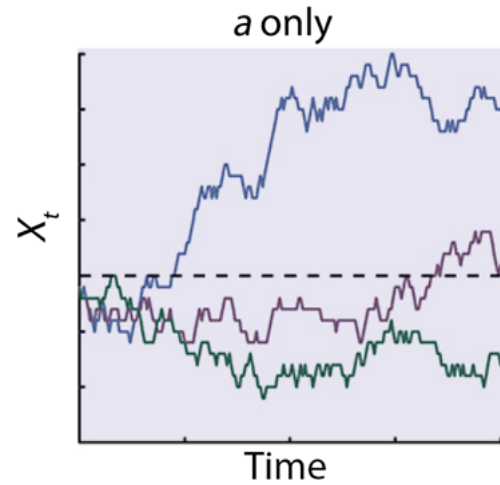
Solution to initial value problem ($u(0, \mathbf{x}) = g(\mathbf{x})$):

Monte Carlo Approximates This Expectation

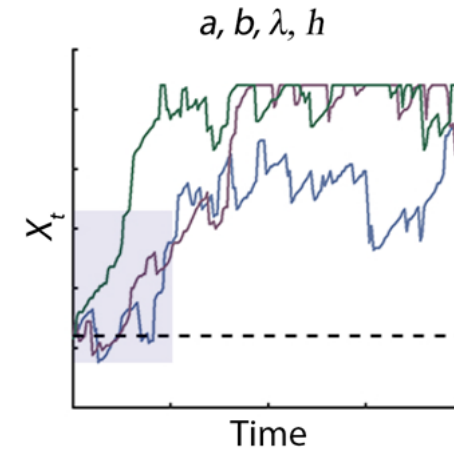
$$u(t, \mathbf{x}) = \mathbb{E} \left[g(\mathbf{X}(t)) \exp \left(\int_0^t c(s, \mathbf{X}(s)) ds \right) + \int_0^t f(s, \mathbf{X}(s)) \exp \left(\int_0^s c(\ell, \mathbf{X}(\ell)) d\ell \right) ds \mid \mathbf{X}(0) = \mathbf{x} \right].$$

Neural MC algorithm can run wide range of stochastic processes

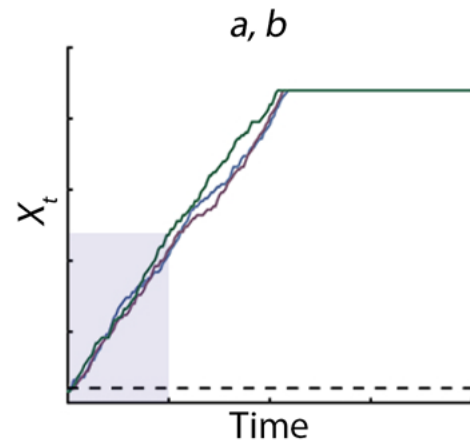
Diffusion



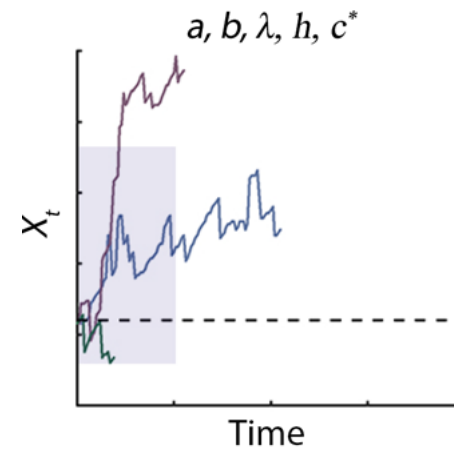
Jump processes



Drift



Absorption / Decay



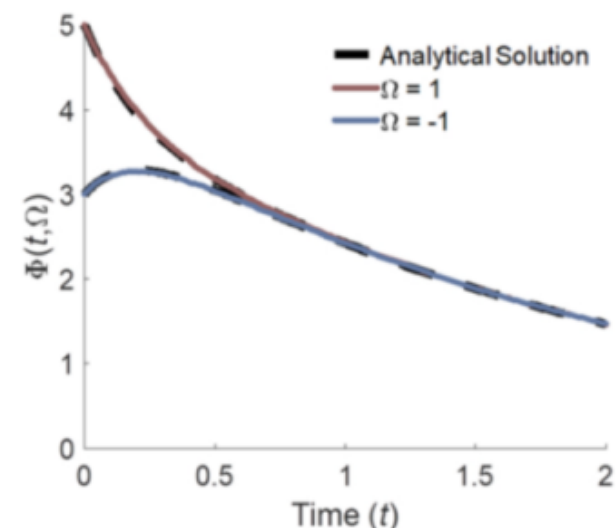
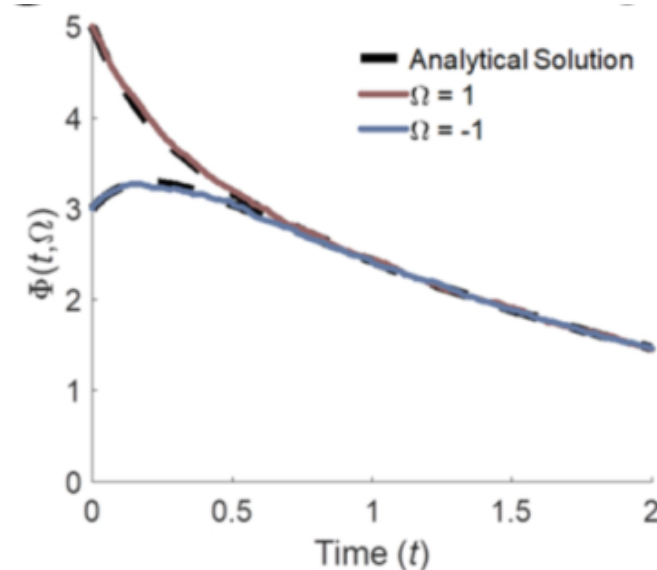
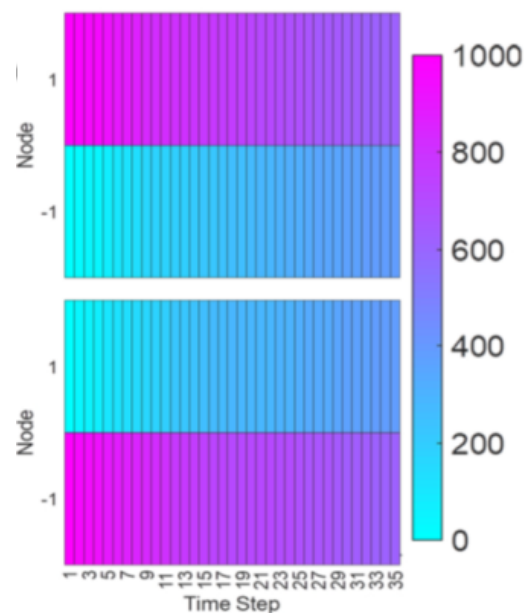
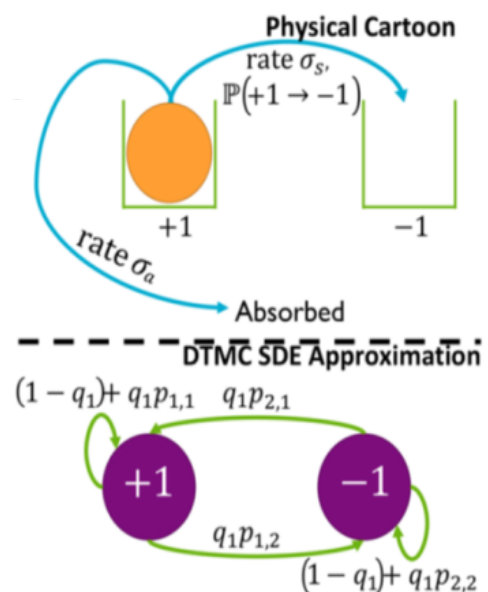
Spiking
Scientific
Computing



Some more applied examples

- Boltzmann state transition
 - Particle can exist in 2 states (+1 or -1) or be absorbed.
 - Implement as simple stochastic process on TrueNorth

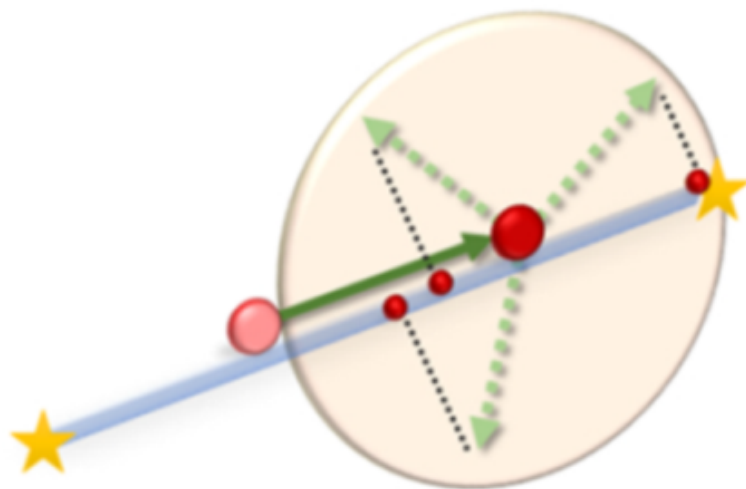
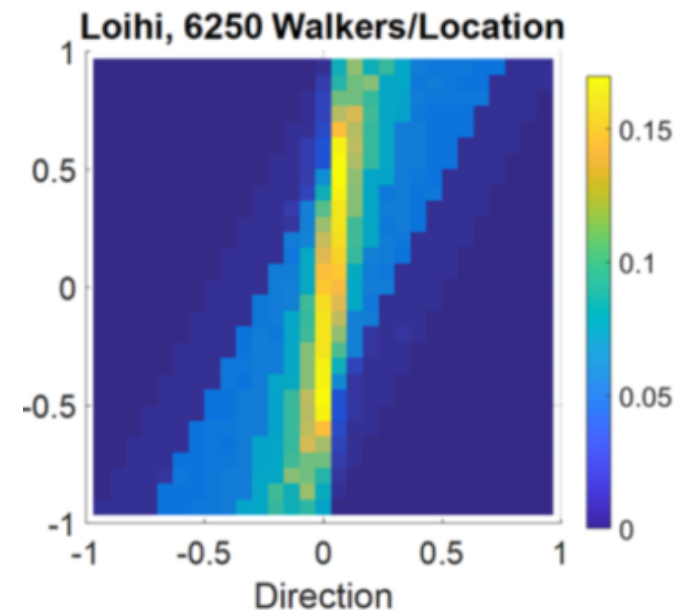
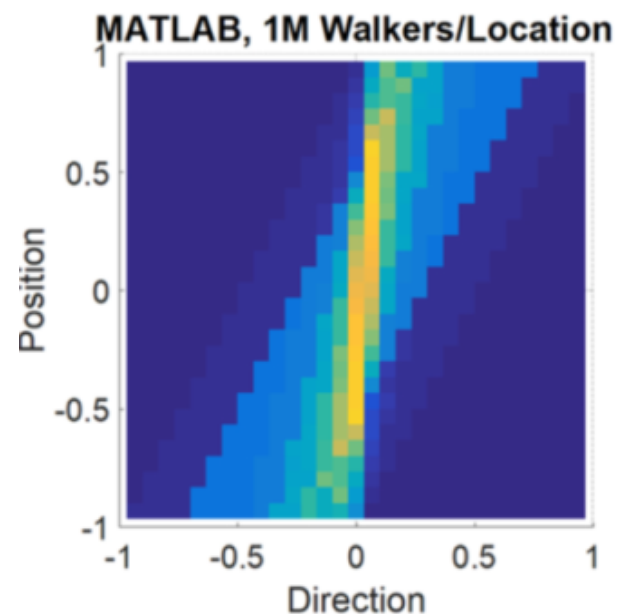
Spiking
Scientific
Computing



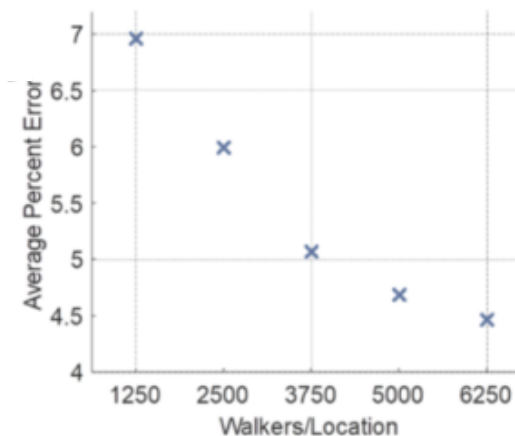


Some more applied examples

- 1D particle transport
 - Particle moves in 2D, only track 1D.
 - At point $x=0$, particle reflects in random direction
 - Track velocity in x-dimension and angle
 - Implemented on Loihi

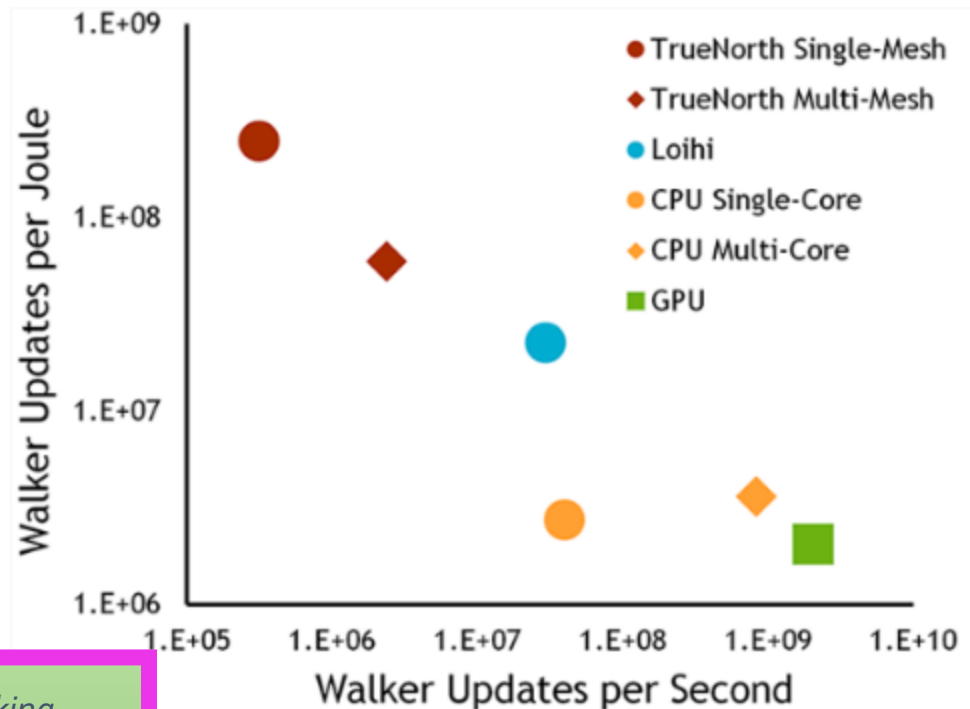


Spiking
Scientific
Computing





Monte Carlo advantage comes from benefits of **spiking** ... still limited by CMOS



Spiking
Scientific
Computing

Broad Applications

Monte Carlo simulations

Discrete Time Markov Chains

Graph Analytics

Graph neural networks

...

Limitations

Stochastic sampling

Neuron and synapse scaling

Configurability of neurons

On chip / off chip
communication



Today's large scale neuromorphic systems are on ***Pareto Frontier*** of computing

Broad class of algorithms fit this tradeoff

- Monte Carlo / Probabilistic
- Graph analytics
- Artificial intelligence
- Optimization

Architectural advantage

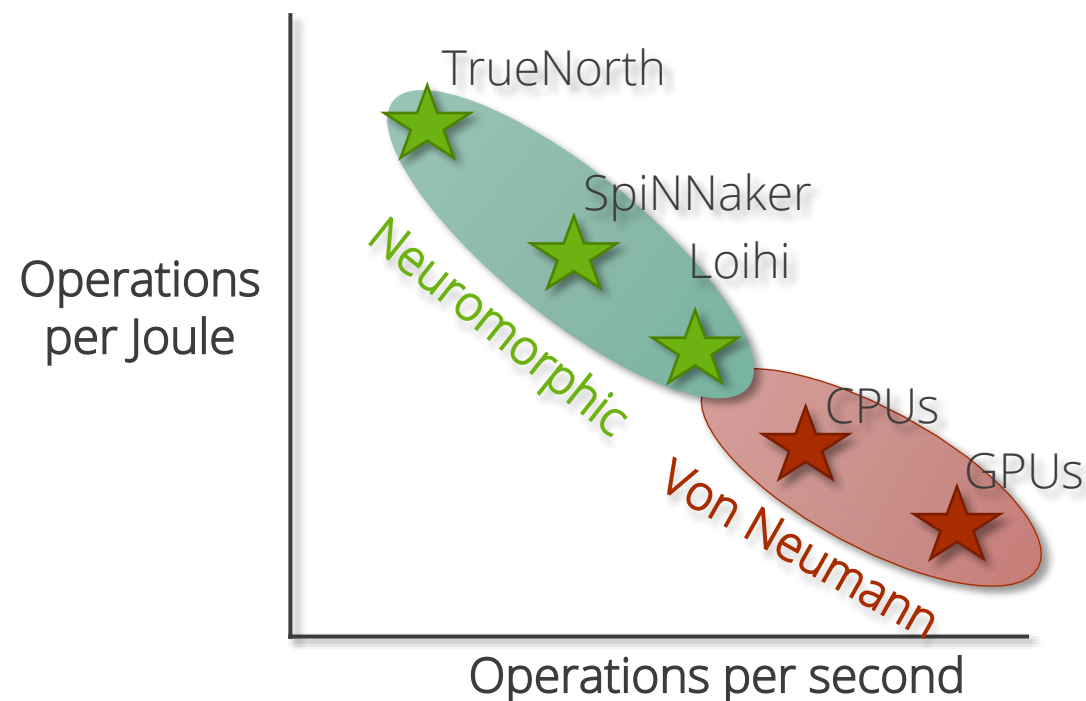
- Event-driven processing
- Massive parallelism

Limitations

- Still CMOS devices
- Architecture is a one time benefit
not an extension to Moore's Law

*Spiking
Scientific
Computing*

If we're honest; who will pick energy efficiency over speed?





Today's large scale neuromorphic systems are on **Pareto Frontier** of computing

Broad class of algorithms fit this tradeoff

- Monte Carlo / Probabilistic
- Graph analytics
- Artificial intelligence
- Optimization

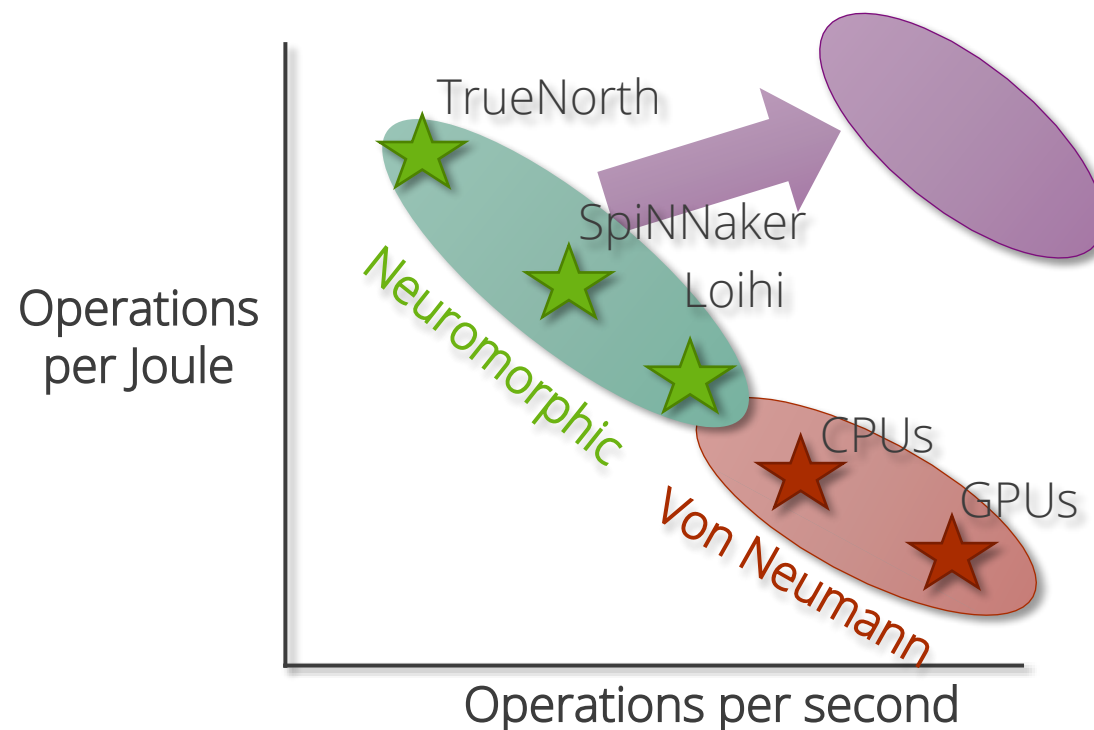
Architectural advantage

- Event-driven processing
- Massive parallelism

Limitations

- Still CMOS devices
- Architecture is a one time benefit
not an extension to Moore's Law

Spiking
Scientific
Computing



Opportunity for Brain-Inspired Materials, Devices & Algorithms
Increasing processing (density, speed, capabilities, etc) while preserving
energy advantage and jump neuromorphic over Pareto Frontier



review articles

DOI:10.1145/3221599

Advances in neurotechnologies are reigniting opportunities to bring neural computation insights into broader computing applications.

BY JAMES B. AIMONE

Neural Algorithms and Computing Beyond Moore's Law

THE IMPENDING DEMISE of Moore's Law has begun to broadly impact the computing research community.³⁸ Moore's Law has driven the computing industry for many decades, with nearly every aspect of society benefiting from the advance of improved computing processors, sensors, and controllers. Behind these products has been a considerable research industry, with billions of dollars invested in fields ranging from computer science to electrical engineering. Fundamentally, however, the exponential growth in computing described by Moore's Law was driven by advances in materials science.^{30,37} From the start, the power of the computer has been limited by the density of transistors. Progressive advances in how to manipulate silicon through advancing lithography methods and new design tools have kept advancing

computing in spite of perceived limitations of the dominant fabrication processes of the time.³⁷

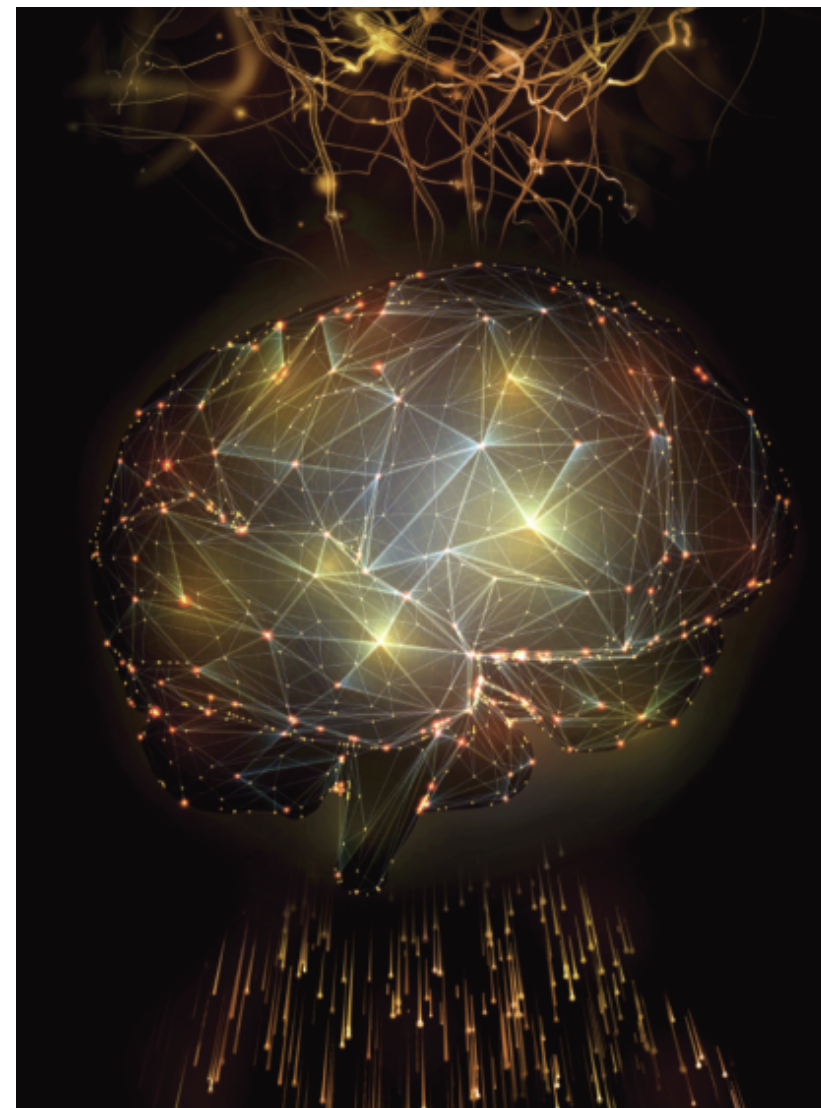
There is strong evidence that this time is indeed different, and Moore's Law is soon to be over for good.^{3,38} Already, Dennard scaling, Moore's Law's lesser known but equally important parallel, appears to have ended.¹¹ Dennard's scaling refers to the property that the reduction of transistor size came with an equivalent reduction of required power.⁴ This has real consequences—even though Moore's Law has continued over the last decade, with feature sizes going from ~65nm to ~10nm; the ability to speed up processors for a constant power cost has stopped. Today's common CPUs are limited to about 4GHz due to heat generation, which is roughly the same as they were 10 years ago. While Moore's Law enables more CPU cores on a chip (and has enabled high power systems such as GPUs to continue advancing), there is increasing appreciation that feature sizes cannot fall much further, with perhaps two or three further generations remaining prior to ending.

Multiple solutions have been presented for technological extension of Moore's Law,^{12,13,14,15} but there are two main challenges that must be addressed. For the first time, it is not immediately evident that future materials

» key insights

- While Moore's Law is slowing down, neuroscience is experiencing a revolution, with technology enabling scientists to have more insights into the brain's behavior than ever before and thus positioning the neuroscience field to provide a long-term source of inspiration for novel computing solutions.
- Extending the reach of brain-inspired computing will not only make current AI methods better, but looking beyond the brain's sensory systems can also expand the reach of AI into new applications.
- Realizing the full potential of brain-inspired computing requires increased collaborations and sharing of knowledge between the neuroscience, computer science, and neuromorphic hardware communities.

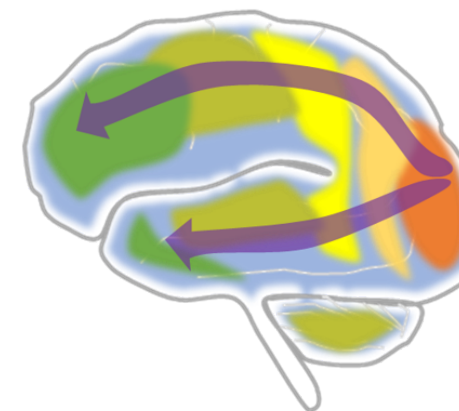
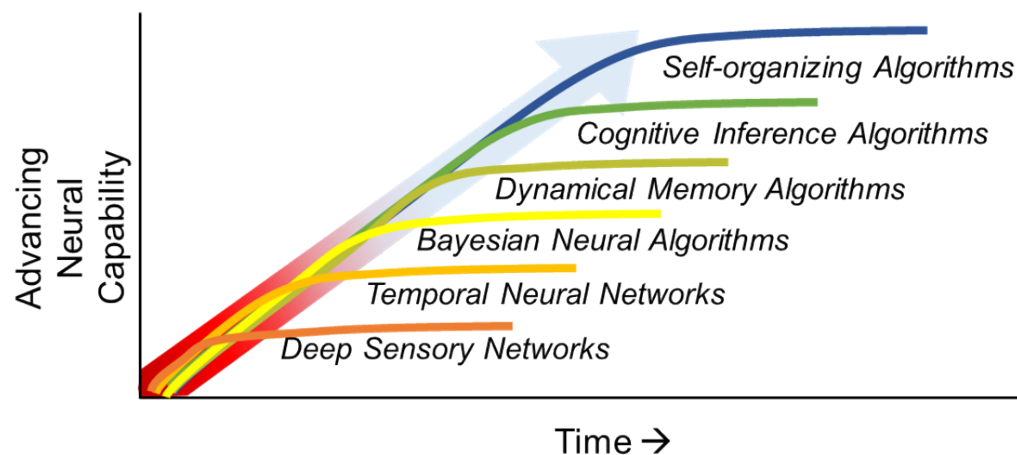
ILLUSTRATION BY NATHAN DAVIS FOR ACM COMMUNICATIONS



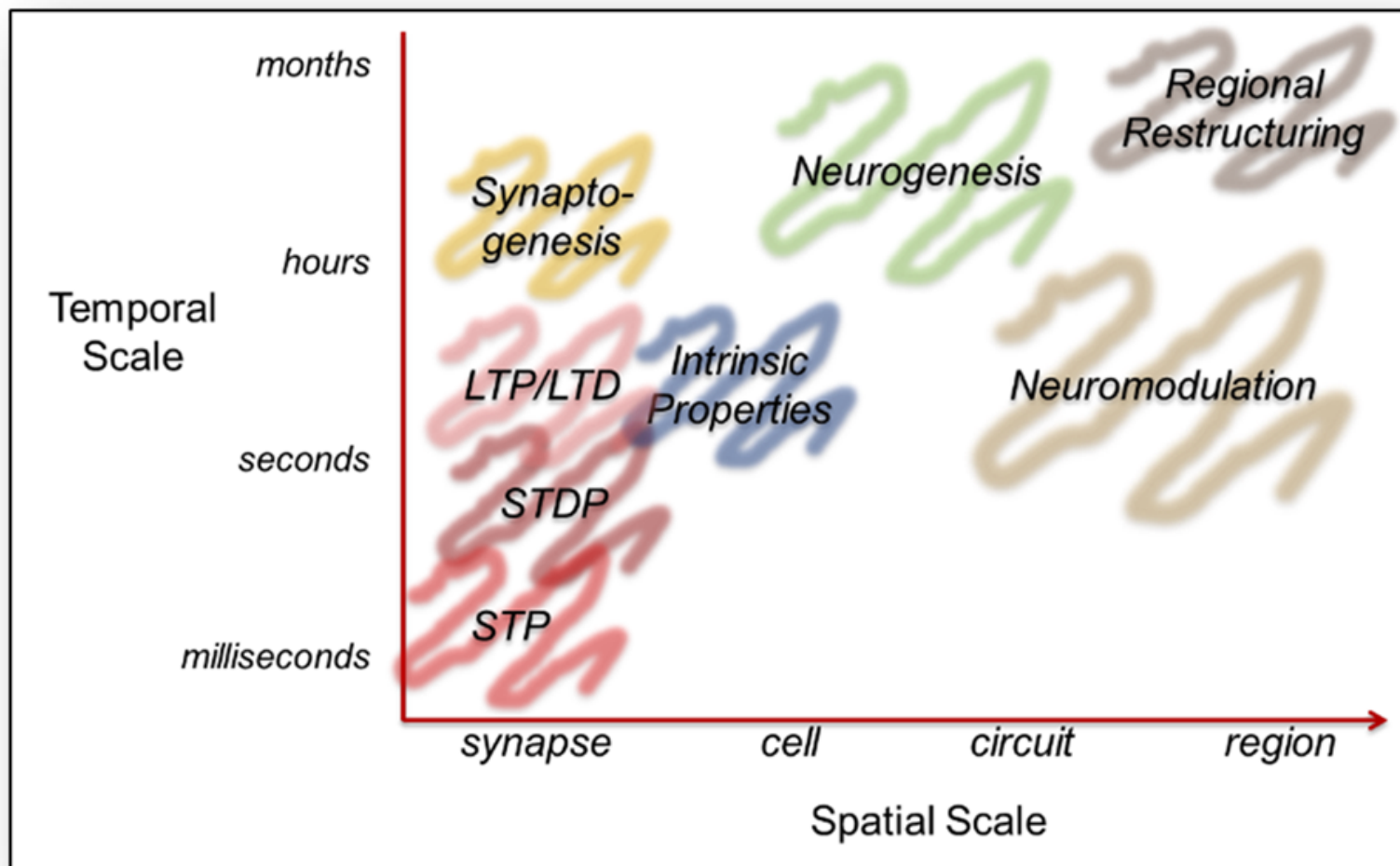
Brain
Inspiration



Algorithm Class	Current Algorithms	Inspiration	Application
Deep Vision Processing	Deep Convolutional Networks (VGG, AlexNet, GoogleNet, etc.), HMax, Neocognitron	Hierarchy of sensory nuclei and early sensory cortices	Static feature extraction (e.g., images) & pattern classification
Temporal Neural Networks	Deep Recurrent Networks (long short-term memory), Hopfield Networks	Local recurrence of most biological neural circuits, especially higher sensory cortices	Dynamic feature extraction (e.g., videos, audio) & classification
Bayesian Neural Algorithms	Predictive Coding, Hierarchical Temporal Memory	Substantial reciprocal feedback between “higher” and “lower” sensory cortices	Inference across spatial and temporal scales
Dynamical Memory and Control Algorithms	Liquid State Machines, Echo State Networks, Neural Engineering Framework	Continual dynamics of hippocampus, cerebellum, and prefrontal and motor cortices	Online learning content-addressable memory & adaptive motor control
Cognitive Inference Algorithms	Reinforcement learning (e.g., Q-learning)	Integration of multiple modalities and memory into prefrontal cortex, which provides top-down influence on sensory processing	Context and experience dependent information processing and decision making
Self-organizing Algorithms	Neurogenesis Deep Learning	Initial development and continuous refinement of neural circuits to specific input and outputs	Automated neural algorithm development for unknown input and output transformations



Brain
Inspiration



Brain
Inspiration



AI systems are slowly starting to incorporate brain-like plasticity

PERSPECTIVE

<https://doi.org/10.1038/s42256-022-00452-0>

nature

machine intelligence

Check for updates

Biological underpinnings for lifelong learning machines

Dhireesha Kudithipudi¹, Mario Aguilar-Simon², Jonathan Babb³, Maxim Bazhenov⁴, Douglas Blackiston^{5,6}, Josh Bongard⁷, Andrew P. Brna⁸, Suraj Chakravarthi Raja⁹, Nick Cheney⁷, Jeff Clune⁸, Anurag Daram¹⁰, Stefano Fusi¹¹, Peter Helfer¹, Leslie Kay¹², Nicholas Ketz¹³, Zsolt Kira¹⁴, Soheil Kolouri¹⁵, Jeffrey L. Krichmar¹⁶, Sam Kriegman¹⁷, Michael Levin¹⁸, Sandeep Madireddy¹⁹, Santosh Manicka⁴, Ali Marjaninejad²⁰, Bruce McNaughton²¹, Risto Miikkulainen²², Zaneta Navratilova²³, Tej Pandit¹, Alice Parker⁴, Praveen K. Pillay²⁴, Sebastian Risi²⁵, Terrence J. Sejnowski^{26,27}, Andrea Soltoggio²⁸, Nicholas Soares²⁹, Andreas S. Tollas³⁰, Dario Urbina-Meléndez³¹, Francisco J. Valero-Cuevas³², Gido M. van de Ven³³, Joshua T. Vogelstein³⁴, Felix Wang³⁵, Ron Weiss³⁶, Angel Yanguas-Gil³⁷, Xinyun Zou³⁸ and Hava Siegelmann³⁹

Biological organisms learn from interactions with their environment throughout their lifetime. For artificial systems to successfully act and adapt in the real world, it is desirable to similarly be able to learn on a continual basis. This challenge is known as lifelong learning, and remains to a large extent unsolved. In this Perspective article, we identify a set of key capabilities that artificial systems will need to achieve lifelong learning. We describe a number of biological mechanisms, both neuronal and non-neuronal, that help explain how organisms solve these challenges, and present examples of biologically inspired models and biologically plausible mechanisms that have been applied to artificial systems in the quest towards development of lifelong learning machines. We discuss opportunities to further our understanding and advance the state of the art in lifelong learning, aiming to bridge the gap between natural and artificial intelligence.

Learning is a defining ability of biological systems, whereby experience leads to behavioural adaptations that improve performance. The past couple of decades have witnessed astonishing advances in the field of machine learning. Nevertheless, a new generation of applications—self-driving cars and trucks, autonomous drones, delivery robots, intelligent handheld and wearable devices, and others that we have not yet imagined—will require a new type of machine intelligence that is able to learn throughout its lifetime. Such machines will need to acquire new skills without compromising old ones, adapt to changes, and apply previously learned knowledge to new tasks—all while conserving limited resources such as computing power, memory and energy. These capabilities are collectively known as lifelong learning (LL).

In contrast to the current generation of intelligent machines, animal species ranging from invertebrates to humans are able to learn continually throughout their lifetime. Neuroscientists and other biologists have proposed several mechanisms to explain this ability, and machine learning researchers have attempted to emulate them in artificial systems, with varying degrees of success. In this Perspective article, we examine our current understanding of how biological organisms learn continually and review the state of the art in biologically inspired LL models. We describe a variety of biological mechanisms, both neuronal and non-neuronal, that can improve our ability to create highly functioning lifelong learning machines. It should be noted that there is also a body of artificial intelligence (AI) research that tackles the lifelong learning problem from a less clearly biological perspective^{1–3}. These can be broadly organized into three types: *rehearsal*, which store or generate data from past tasks for replay^{4–6}; *architectural*, which expand the model parameters^{7–11}; and *regularization-based* approaches, which penalize changes to parameters important to past tasks^{12–15} or use meta-learning¹⁶. Such models, which are not directly inspired by a biological mechanism, fall outside the scope of this Perspective.

In this Perspective, we will (1) identify a set of key features of lifelong learning; (2) provide an overview of biological mechanisms that are believed to be involved in realizing these features; and (3) review research in which analogous mechanisms have been implemented in machine learning models with the aim of realizing lifelong learning capabilities in artificial systems. We conclude with a look at future challenges and opportunities.

¹University of Texas at San Antonio, San Antonio, TX, USA. ²Intelligent Systems Laboratory, Teledyne Scientific, RTP, NC, USA. ³Massachusetts Institute of Technology, Boston, MA, USA. ⁴University of California at San Diego, La Jolla, CA, USA. ⁵Allen Discovery Center, Tufts University, Medford, MA, USA. ⁶Wyss Institute, Harvard University, Cambridge, MA, USA. ⁷University of Vermont, Burlington, VT, USA. ⁸University of Southern California, Los Angeles, CA, USA. ⁹University of British Columbia, Vancouver, BC, Canada. ¹⁰Columbia University, New York, NY, USA. ¹¹University of Chicago, Chicago, IL, USA. ¹²HRG, Laboratoire, Maribou, CA, USA. ¹³Georgia Institute of Technology, Atlanta, GA, USA. ¹⁴Vanderbilt University, Nashville, TN, USA. ¹⁵University of California, Irvine, CA, USA. ¹⁶Argonne National Laboratory, Lemont, IL, USA. ¹⁷The University of Texas at Austin, Austin, TX, USA. ¹⁸MIT, University of Copenhagen, Copenhagen, Denmark. ¹⁹Salk Institute for Biological Studies, La Jolla, CA, USA. ²⁰Loughborough University, Loughborough, UK. ²¹Rochester Institute of Technology, Rochester, NY, USA. ²²Baylor College of Medicine, Houston, TX, USA. ²³Johns Hopkins University, Baltimore, MD, USA. ²⁴Sandia National Laboratories, Albuquerque, NM, USA. ²⁵University of Massachusetts, Amherst, MA, USA. ²⁶Re-mail: dr.sejnowski@utoronto.ca

NATURE MACHINE INTELLIGENCE | VOL. 4 | MARCH 2022 | 196–203 | www.nature.com/nature-machine-intelligence

Biologically inspired mechanisms

Key features					
Transfer and adaptation	Overcoming catastrophic forgetting	Exploiting task similarity	Task-agnostic learning	Noise tolerance	Resource efficiency and sustainability
Neurogenesis	169–174	234	161		174,201,202
Episodic replay	54,175,176,179,180		54,176	176,177	53,54,175,176,179,180,203
Metaplasticity	67,89,181–185		7,181,185,198		89,181–183,198
Neuromodulation	70,78,84–86,88,89,157,159,160	78,79,84,89,164	89	78,159	78,158,199
Context-dependent perception and gating	78,79,158,161–167	78,168	79,162–166	70,161	158,162,163
Hierarchical distributed systems			188–191		113,191,200
Cognition outside the brain			195–197		
Reconfigurable organisms	139		139,147		139,147
Multisensory integration			152,155,192,193		113,162

Brain Inspiration



A concrete future direction: Brain-inspired systems may need to embrace stochasticity

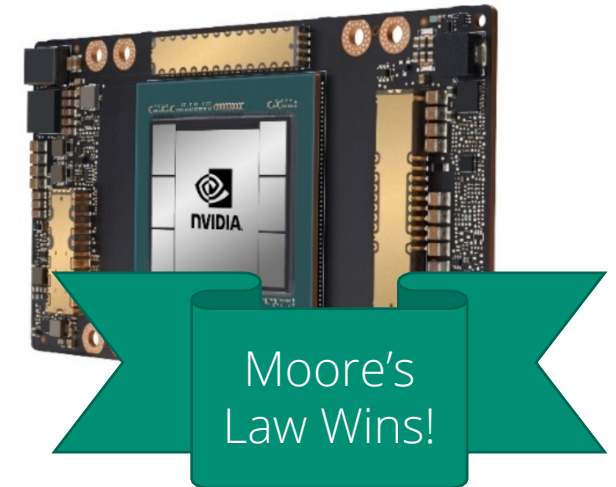


Has the tremendous success of deterministic computing left probabilistic applications behind?

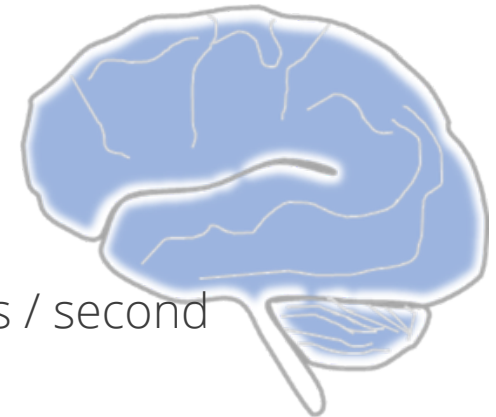
Which approach is best to interpret a clear input?



~400 W
~ 10^{13} - 10^{14} FLOPS
Fully deterministic



~20 W
~ 10^{15} synaptic events / second
Fully stochastic



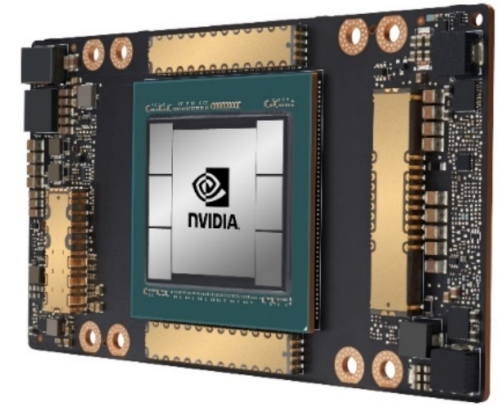


Has the tremendous success of deterministic computing left probabilistic applications behind?

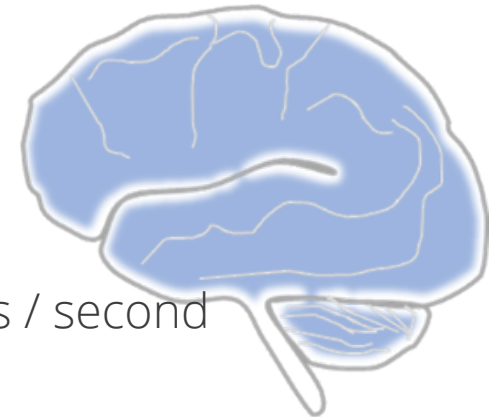
Which approach is best to interpret an ambiguous input?



~400 W
~ 10^{13} - 10^{14} FLOPS
Fully deterministic



~20 W
~ 10^{15} synaptic events / second
Fully stochastic



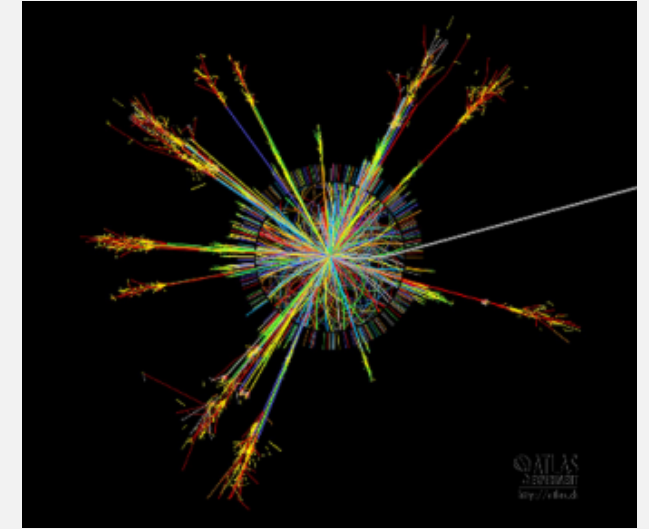
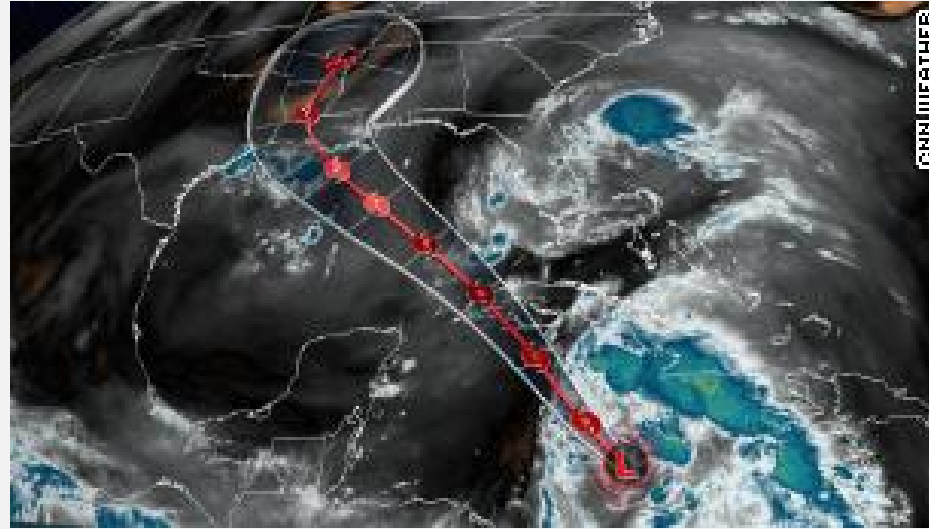


Computing applications face challenges in uncertainty



Artificial Intelligence

- Bayesian neural networks are appealing yet often computationally intractable

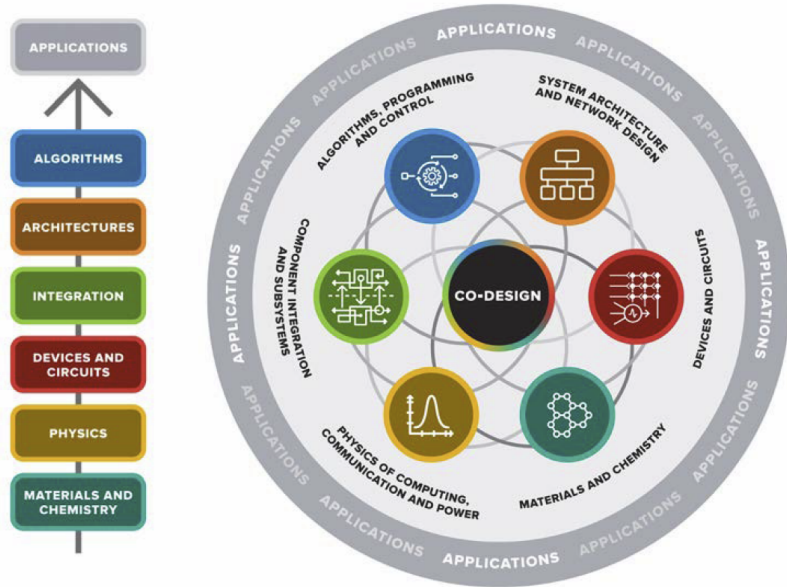


Modeling and Simulation

- Modeling uncertainties is critical in the use of even fully deterministic simulations
- Many applications are inherently stochastic in their physics and are best modeled using probabilistic methods

CO-designed Improved Neural Foundations Leveraging Inherent Physics Stochasticity (COINFLIPS)

- Office of Science Co-Design in Microelectronics program
 - Co-funded through ASCR and BES, participation by NP, HEP, and FES



To enable new generations of energy-efficient computing systems over the next decade, a complete reconceptualization of the science and technology underlying the microelectronics co-design approach is needed to integrate emerging devices, materials, interconnects, and non-linear phenomena with the needs of scientific computing applications.



U.S. DEPARTMENT OF
ENERGY

Office of
Science



Sandia
National
Laboratories



OAK RIDGE
National Laboratory





CO-designed Improved Neural Foundations Leveraging Inherent Physics Stochasticity (COINFLIPS)

- Office of Science Co-Design in Microelectronics program
 - Co-funded through ASCR and BES, participation by NP, HEP, and FES
- ~COINFLIPS is partnering with a growing number of organizations
 - Andy Kent @ New York University
 - Jean Anne Incorvia @ University of Texas Austin
 - Katie Schuman @ University of Tennessee
 - Prasanna Date @ Oak Ridge National Laboratory
 - Les Bland @ Temple University



U.S. DEPARTMENT OF
ENERGY

Office of
Science



Sandia
National
Laboratories



OAK RIDGE
National Laboratory





We are benefitting from 70 years of microelectronics that embrace ***deterministic*** components to solve ***deterministic*** problems

COINFLIPS sees an opportunity to embrace ***stochastic*** computing to solve ***uncertainty*** problems



Today's computers emulate uncertainty by using pseudo-random number generation



“Any one who considers arithmetical methods of producing random digits is, of course, in a state of sin.”

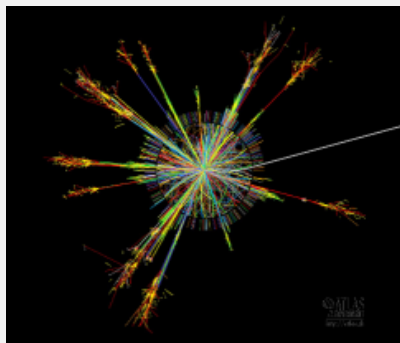
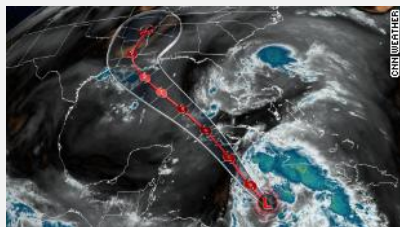
John von Neumann, 1951

70 years later...

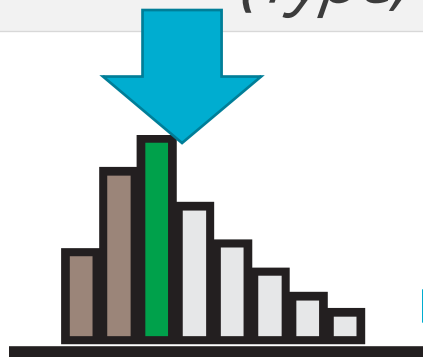
- Pseudo-RNGs can be quite effective, and do offer some advantages in verification, etc.
- But they are expensive, and when they go wrong the implications can be disastrous



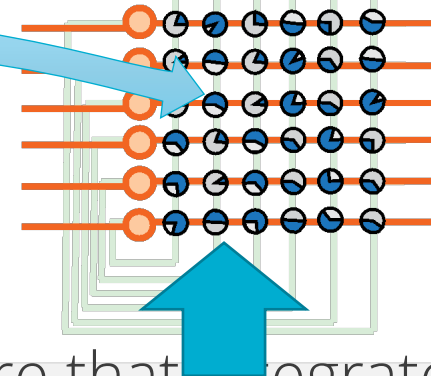
COINFLIPS aims to integrate true random number generators using stochastic devices into neuromorphic architectures



Improved Random Number Generation
(*Type, Quantity, Quality*)



And sample that number *where* it is needed within the computation



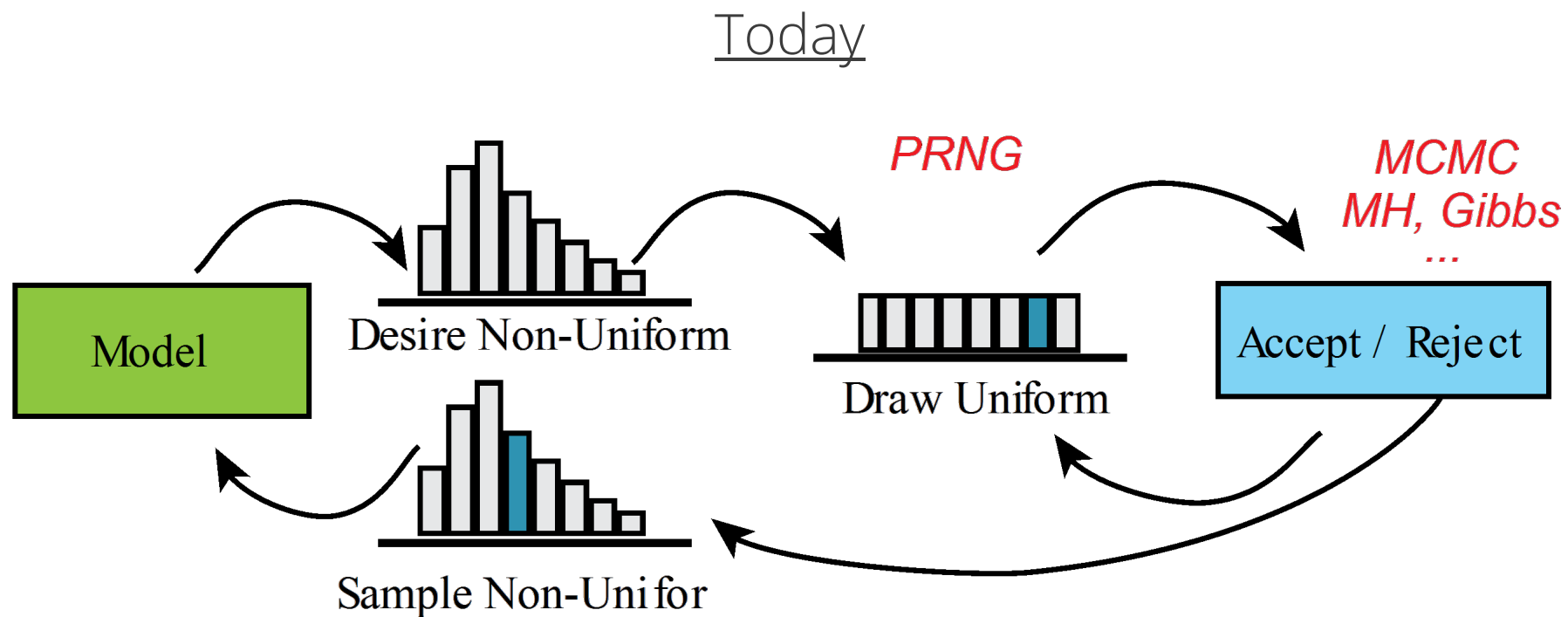
Sample a random number from the *exact* distribution we require

Neuromorphic architecture that integrates ubiquitous stochastic devices with computing and memory

COINFLIPS aims to improve both speed and energy of probabilistic computing applications



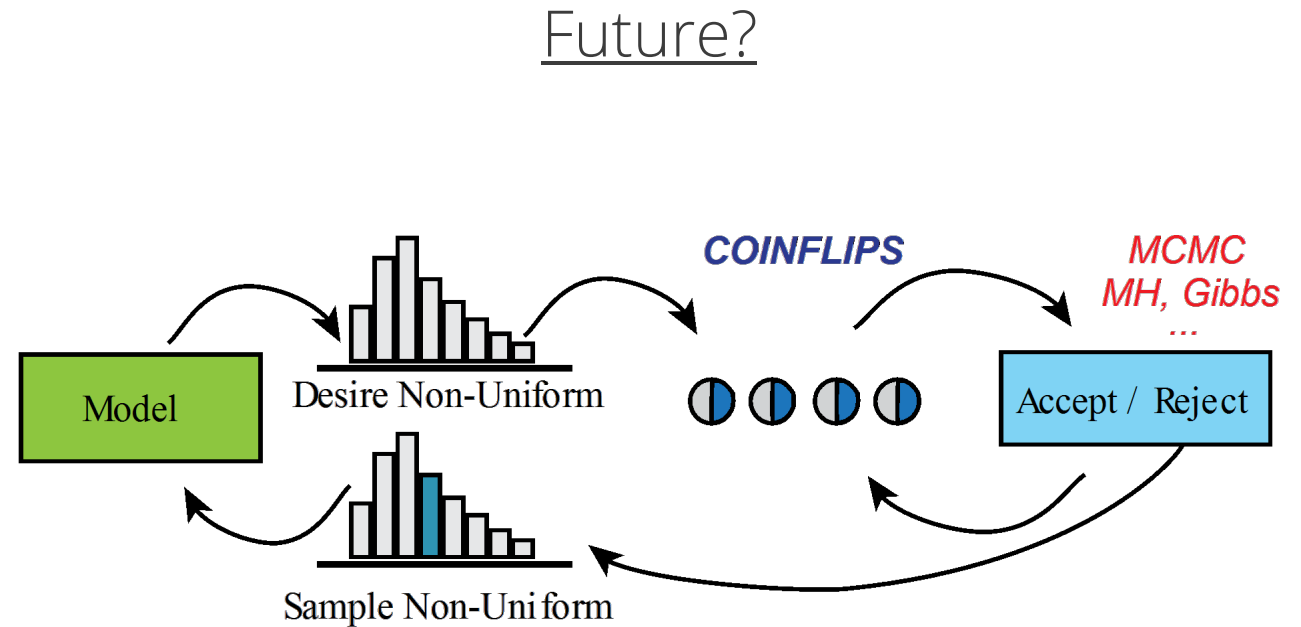
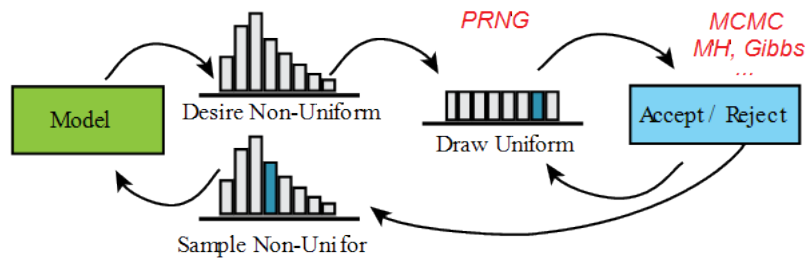
Evaluate opportunity of a probabilistic computing paradigm



COINFLIPS

Evaluate opportunity of a probabilistic computing paradigm

Today

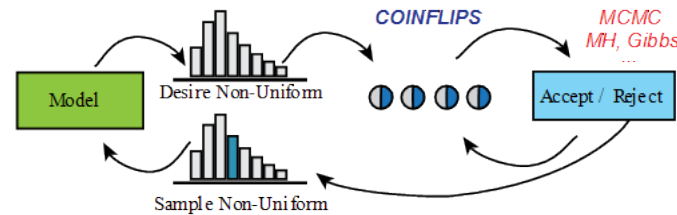


Step 1: Draw suitable uniform RNs from hardware

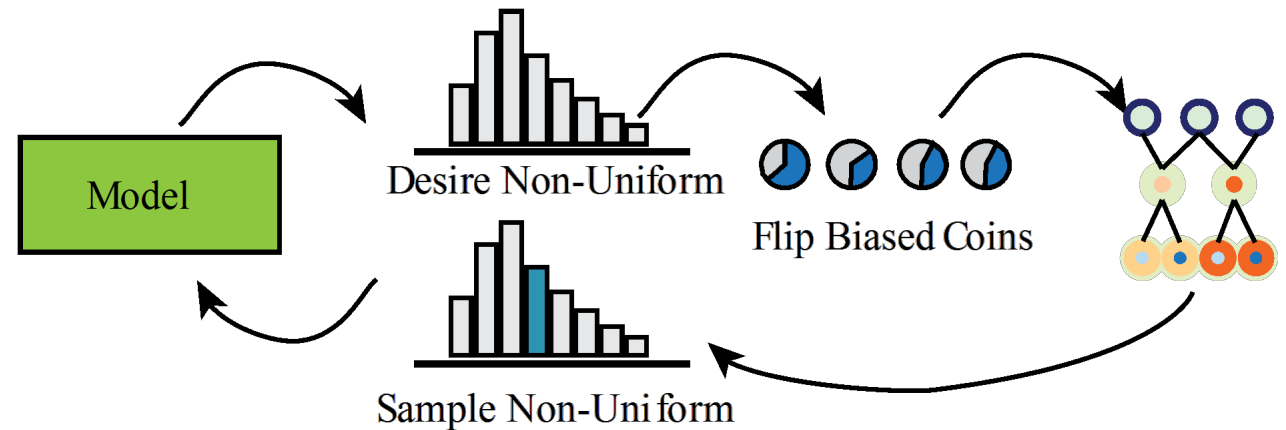
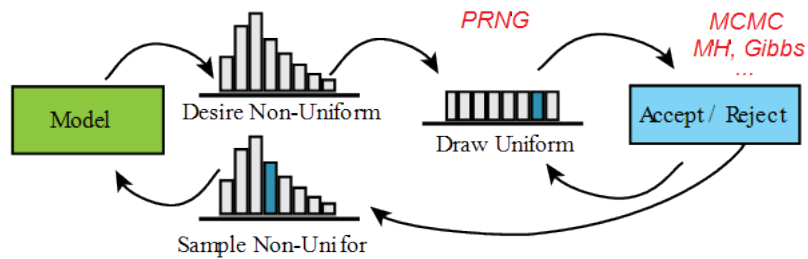
COINFLIPS

Evaluate opportunity of a probabilistic computing paradigm

Future?



Today



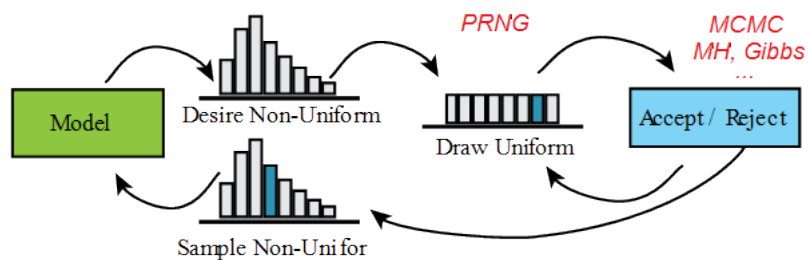
Step 2: Draw suitable model-specific RNs from hardware

COINFLIPS



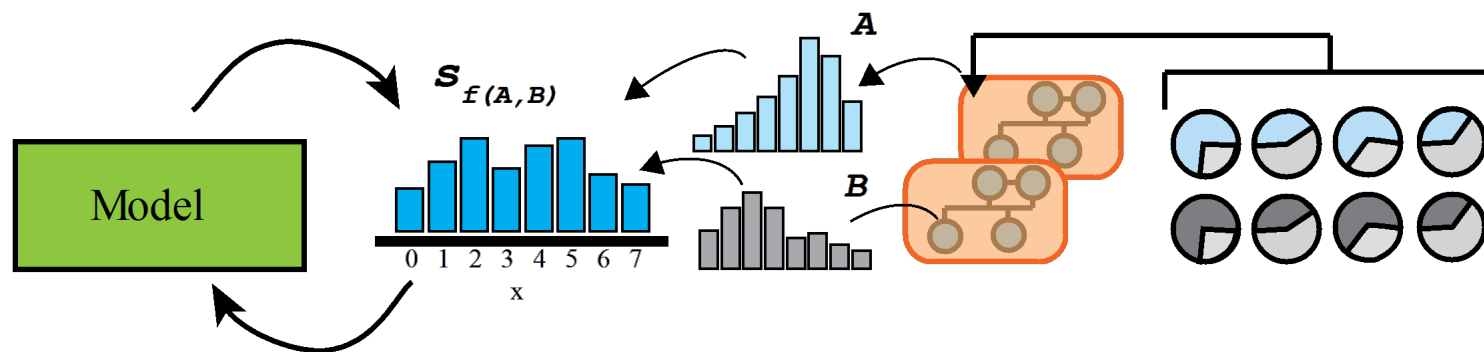
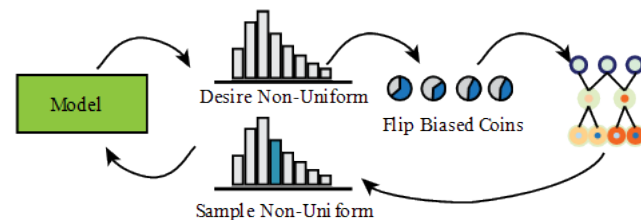
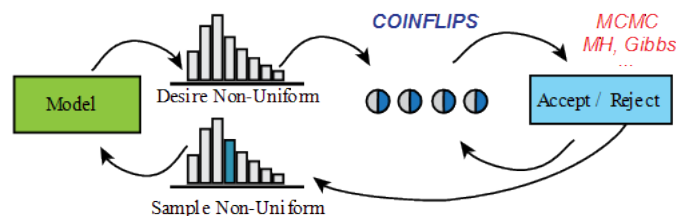
Evaluate opportunity of a probabilistic computing paradigm

Today



COINFLIPS

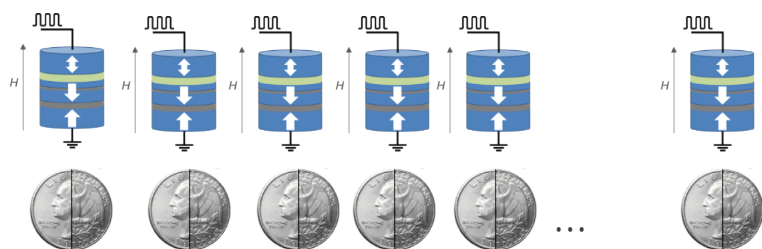
Future?



Step 3: Integrate hardware-enabled random sampling into computation

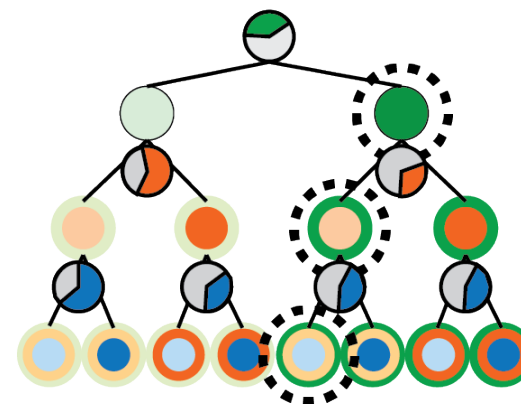
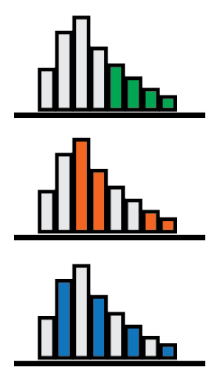


Mapping Coinflips to Arbitrary Distributions



Many devices
flipping at one
time

Naively, we can expand a binary tree with probabilities
to describe any distribution

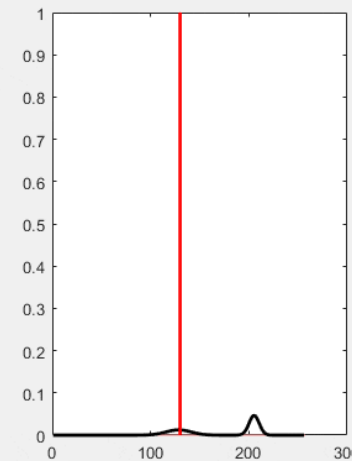
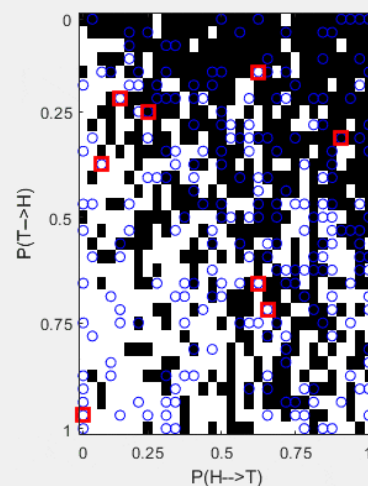
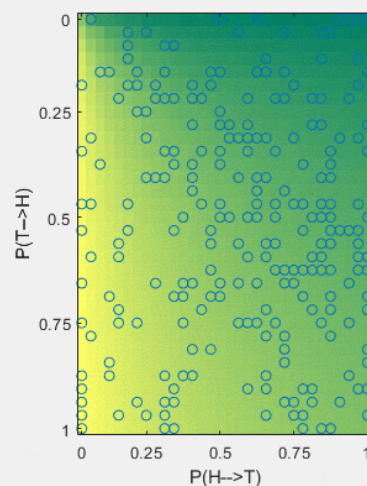


$x = 1$ — —

$x = 1 0$ —

$x = 1 0 0$

Probabilistic
Neural Theory
and Algorithms

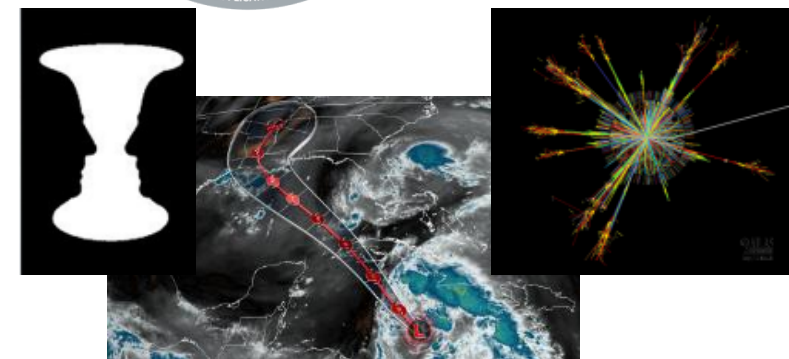
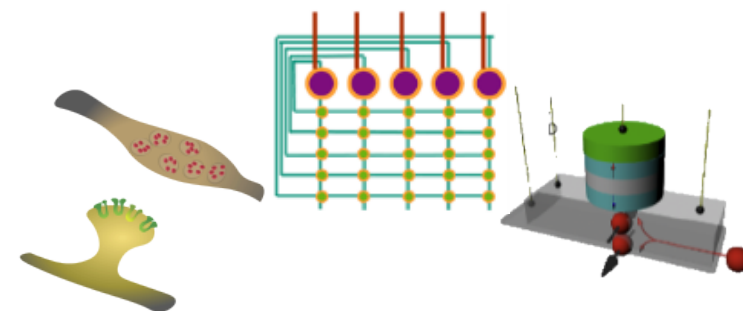


Simulation...



Summary: Probabilistic computing is perhaps an ideal target for exploring potential for microelectronics co-design

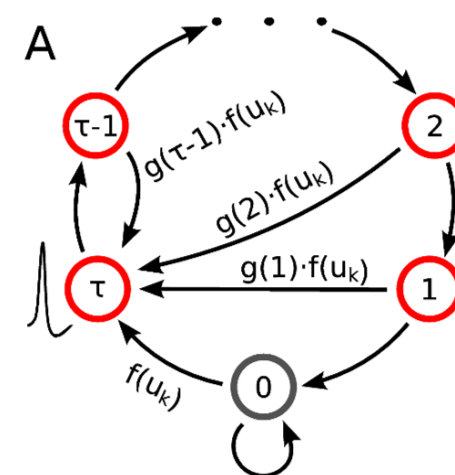
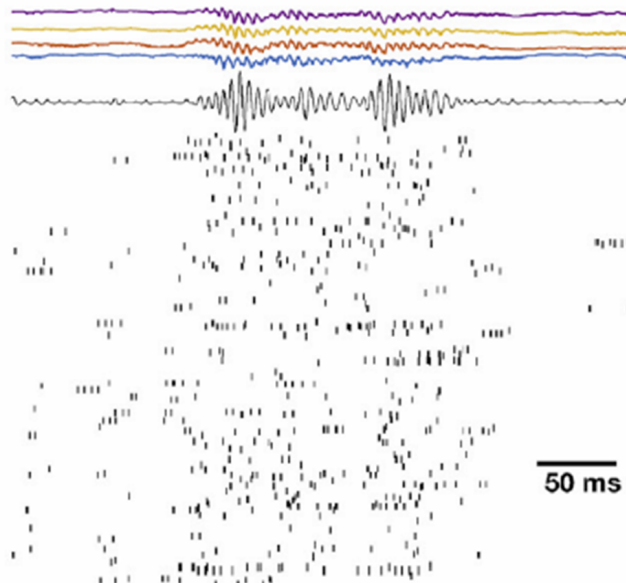
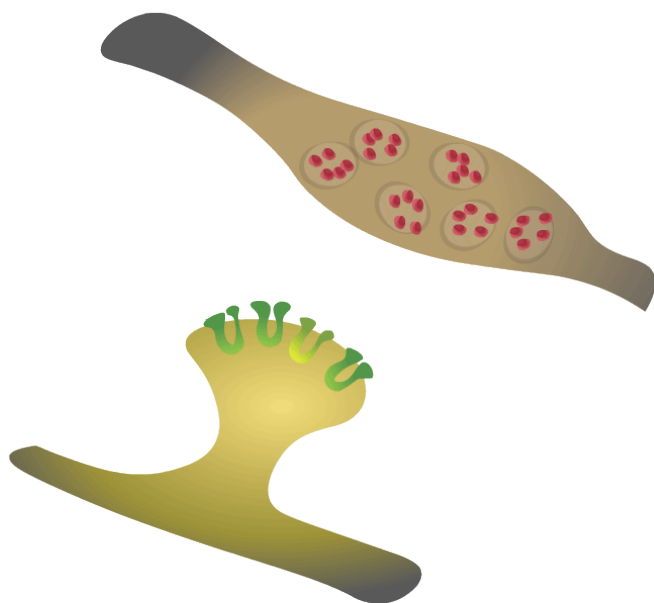
- All aspects of microelectronics (from materials to applications) have something to contribute
 - *Can show benefits from innovation at all scales*
- Stochastic devices
+ neuromorphic parallelism
= broad application impact
 - *Both Mod-Sim and AI stand to benefit*
- Opportunity to consider important aspects of computing up front
 - *Address issues such as I/O, programmability, and theory from the onset, as opposed to after-the-fact*





Backup: Neuroscience and stochastic computation

There is a long history of viewing neuroscience through a stochastic computation perspective. Most of this history is independent of envisioned computing applications



"Independent sources of quantal variability at single glutamatergic synapses" Franks KM, Stevens CF, Sejnowski TJ. *Neurosci.* 2003

"Hippocampal Reactivation of Random Trajectories Resembling Brownian Diffusion" Stella F et al., *Neuron.* 2017

"Neural Dynamics as Sampling: A model of stochastic computation ..." Buesing L et al., *PLOS Computational Biology.* 2011



Has the tremendous success of deterministic computing left probabilistic applications behind?

Stochasticity reveals contrast in computing approaches

- Modern microelectronics spends tremendous resources in enforcing determinism
- The brain embraces and controls stochasticity across spatial and time scales

Developing probabilistic computing to address probabilistic applications

- **COINFLIPS** is combining stochastic devices with neuromorphic architectures
- Co-design is proving invaluable in developing this novel paradigm for microelectronics

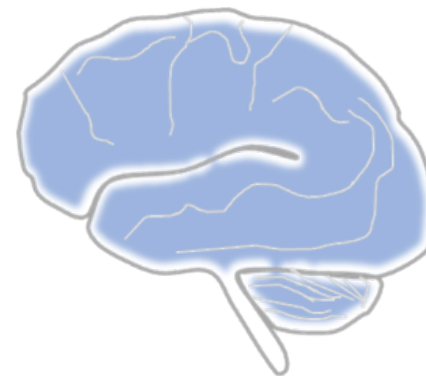
Which approach is best to interpret an ambiguous input?



~20 W

~ 10^{15} synaptic events / second

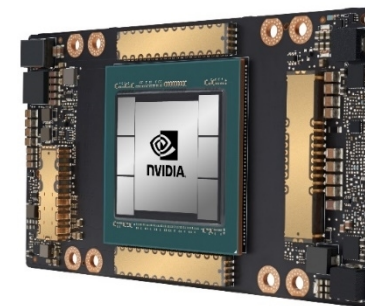
Fully stochastic



~400 W

~ 10^{13} - 10^{14} FLOPS

Fully deterministic



Thanks

